

# Computer Vision

Pietro Perona  
California Institute of Technology

DS@HEP 2017  
Fermilab - 8 May 2017



What?

Where?



# Scene understanding



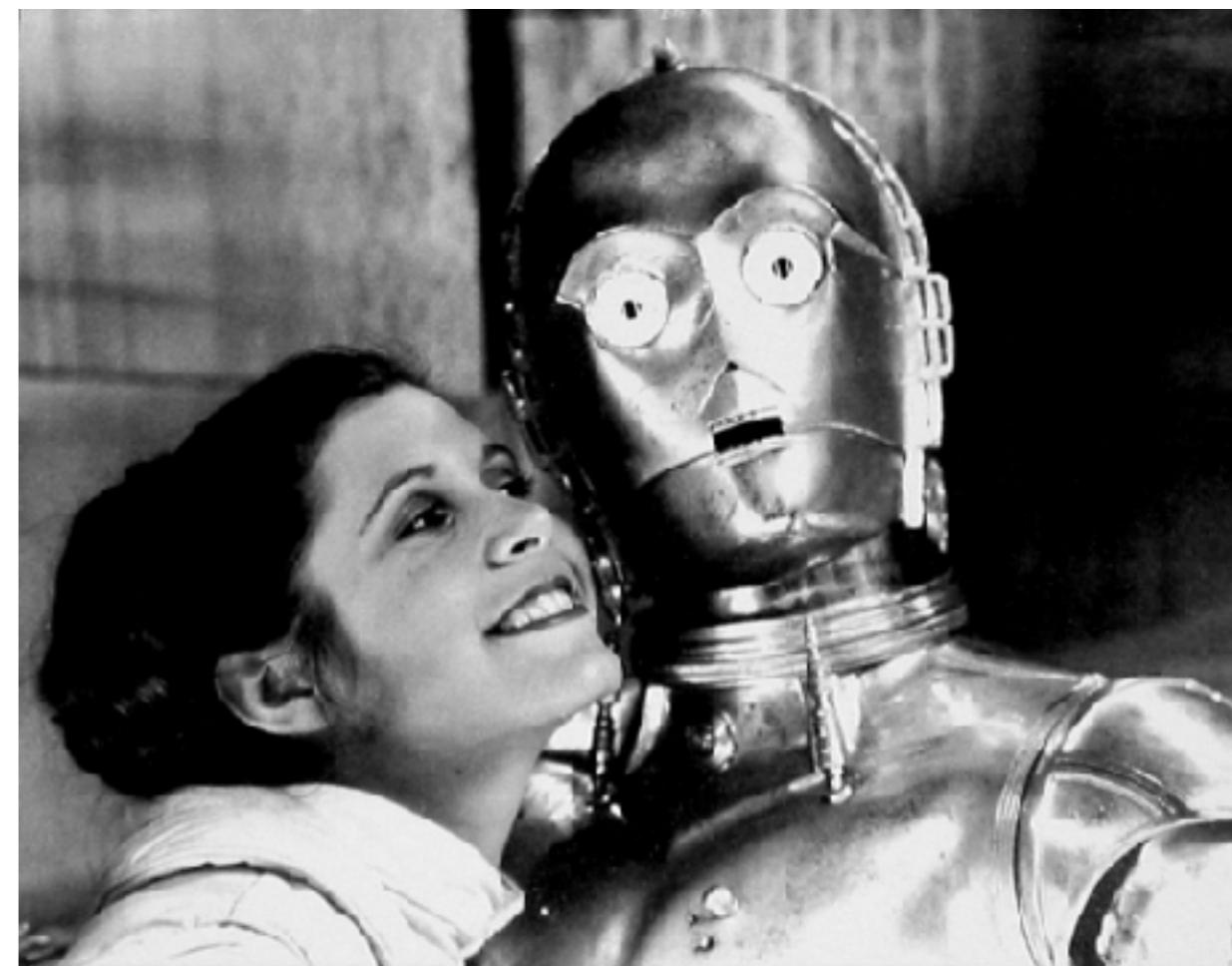
# Scene understanding

On the Italian Alps. Foggy but not too cold. Light breeze. Boy in danger of falling. Parent must be nearby.

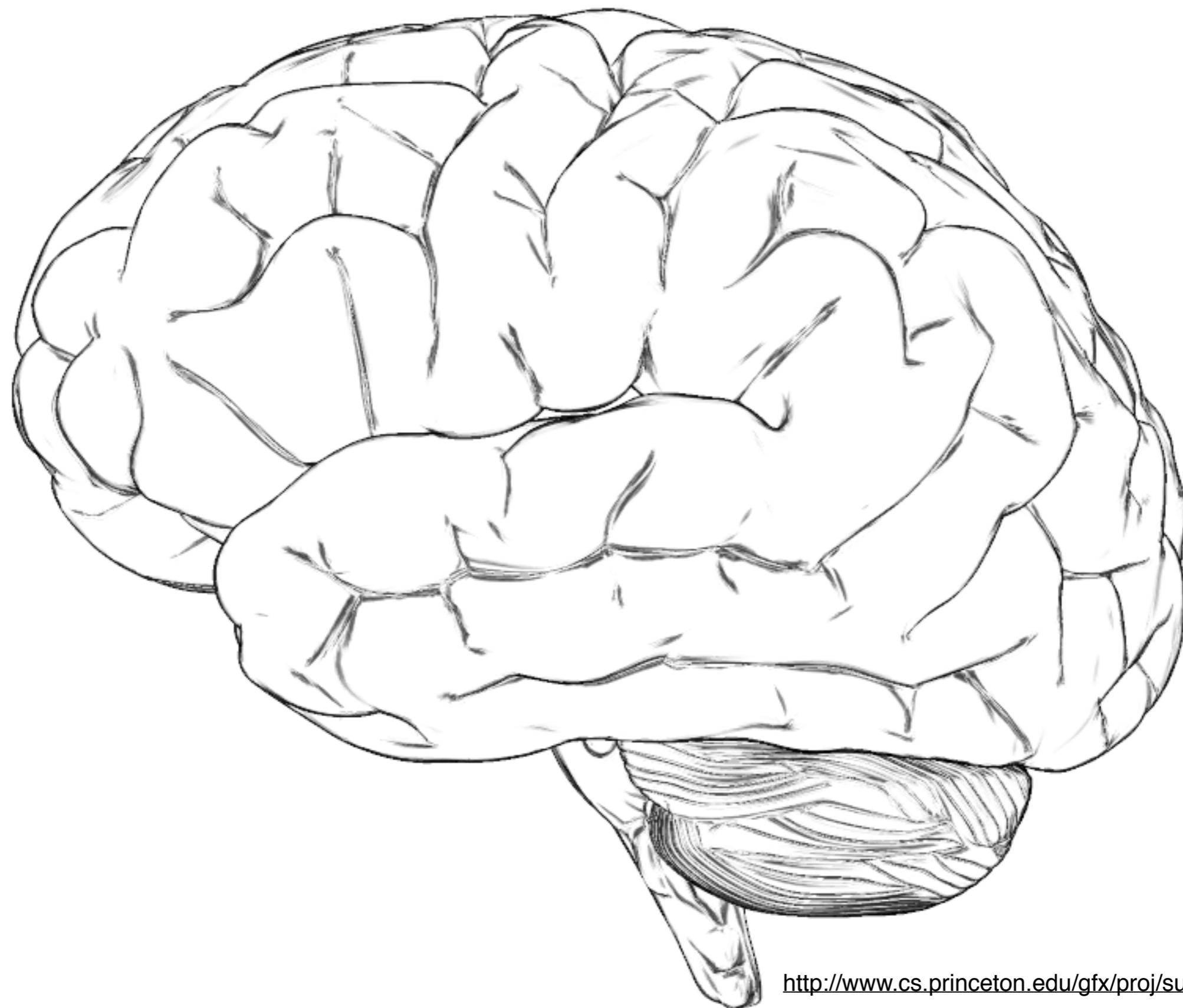


Why study vision?

# Many applications

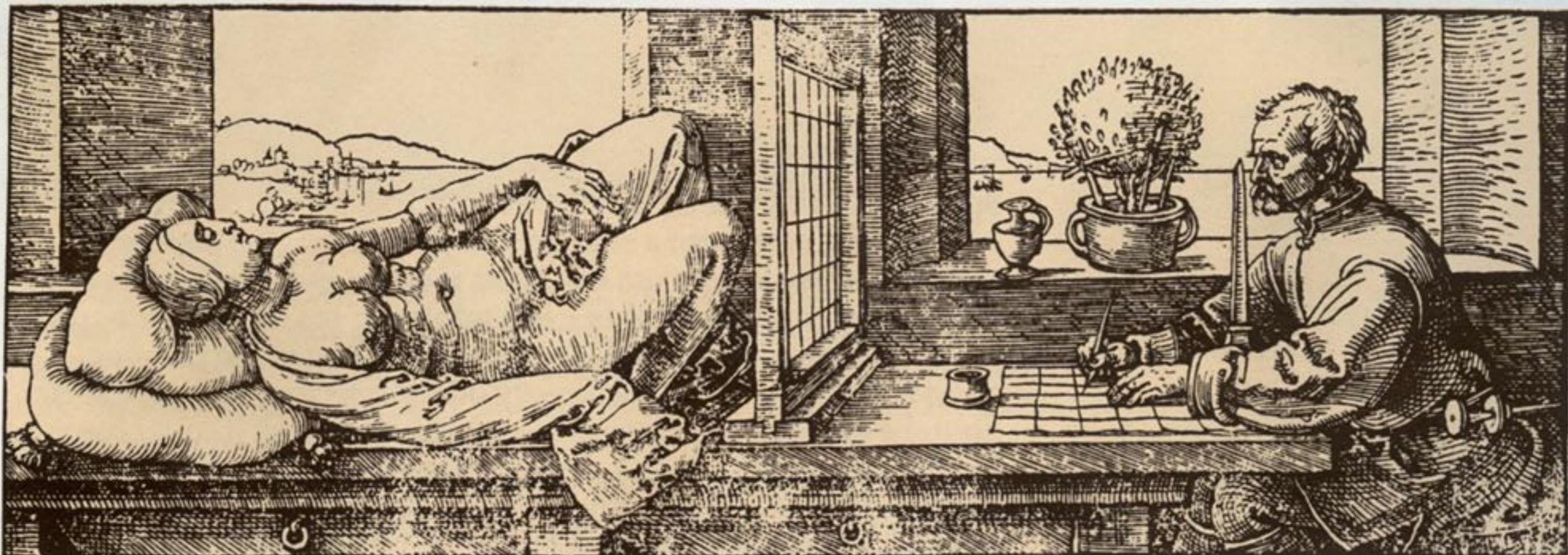


# Understanding the brain



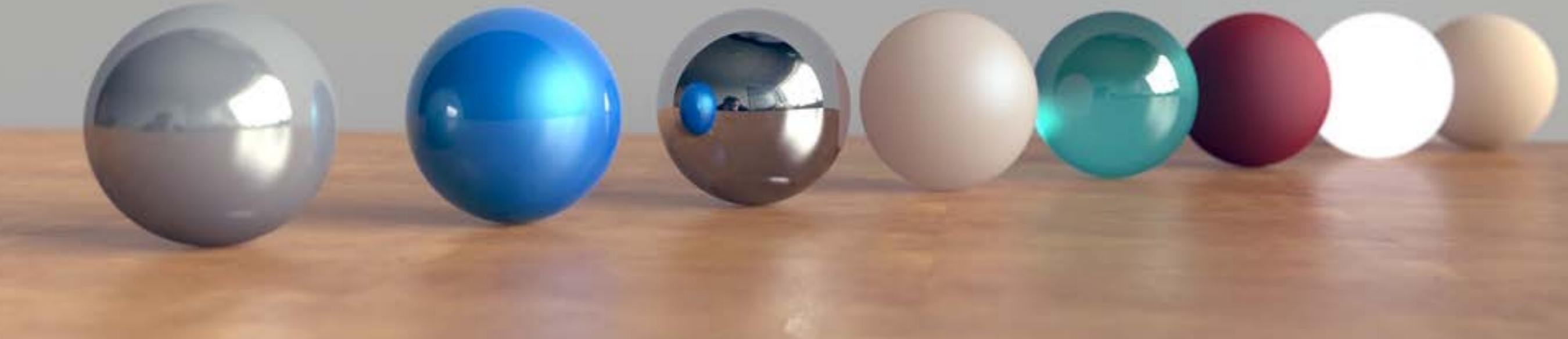
Vision as an inverse problem

# geometry of image formation

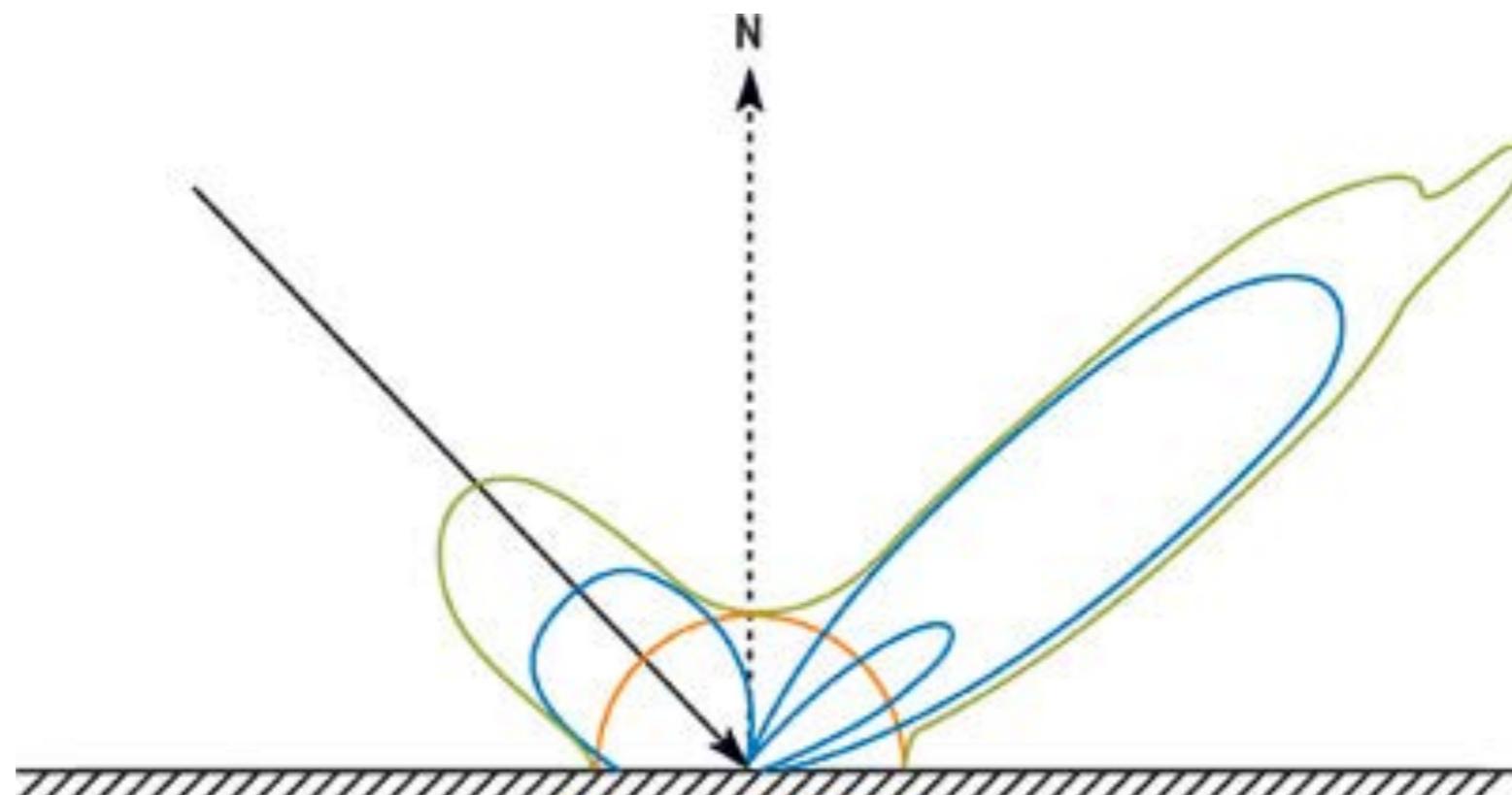


[A. Dürer 1525]

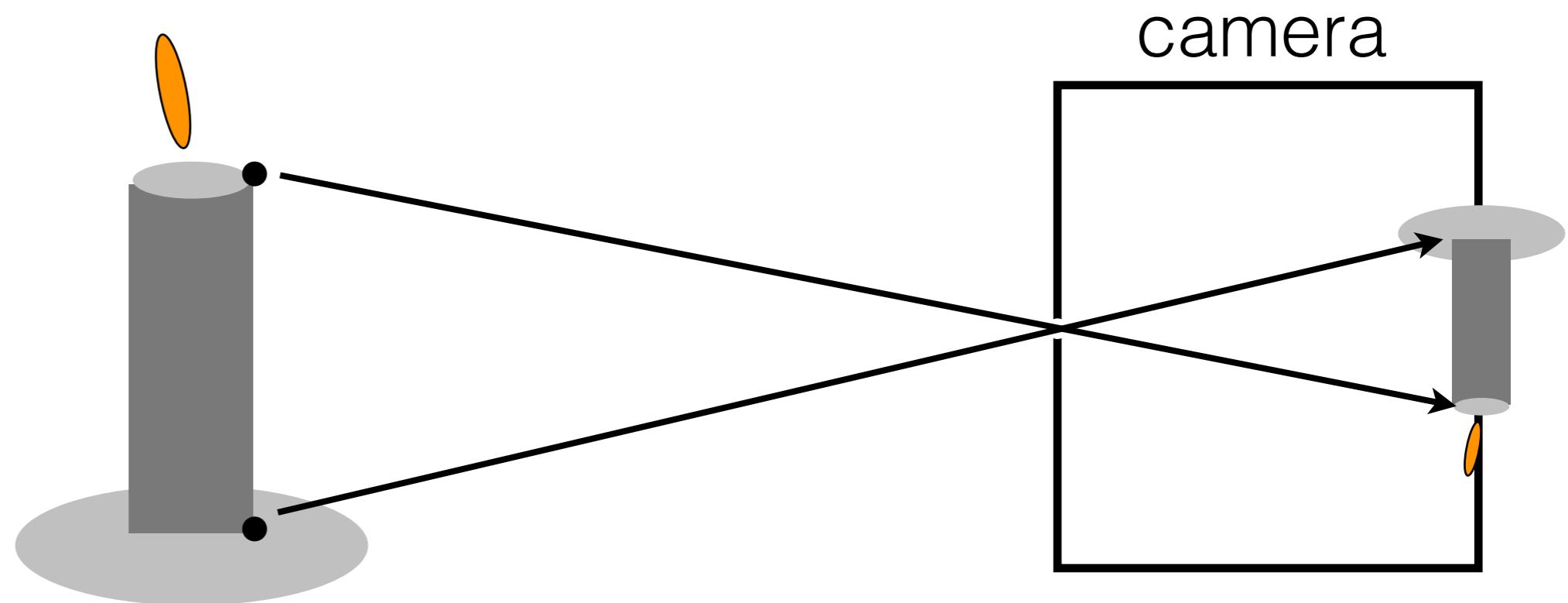
# photometry of image formation



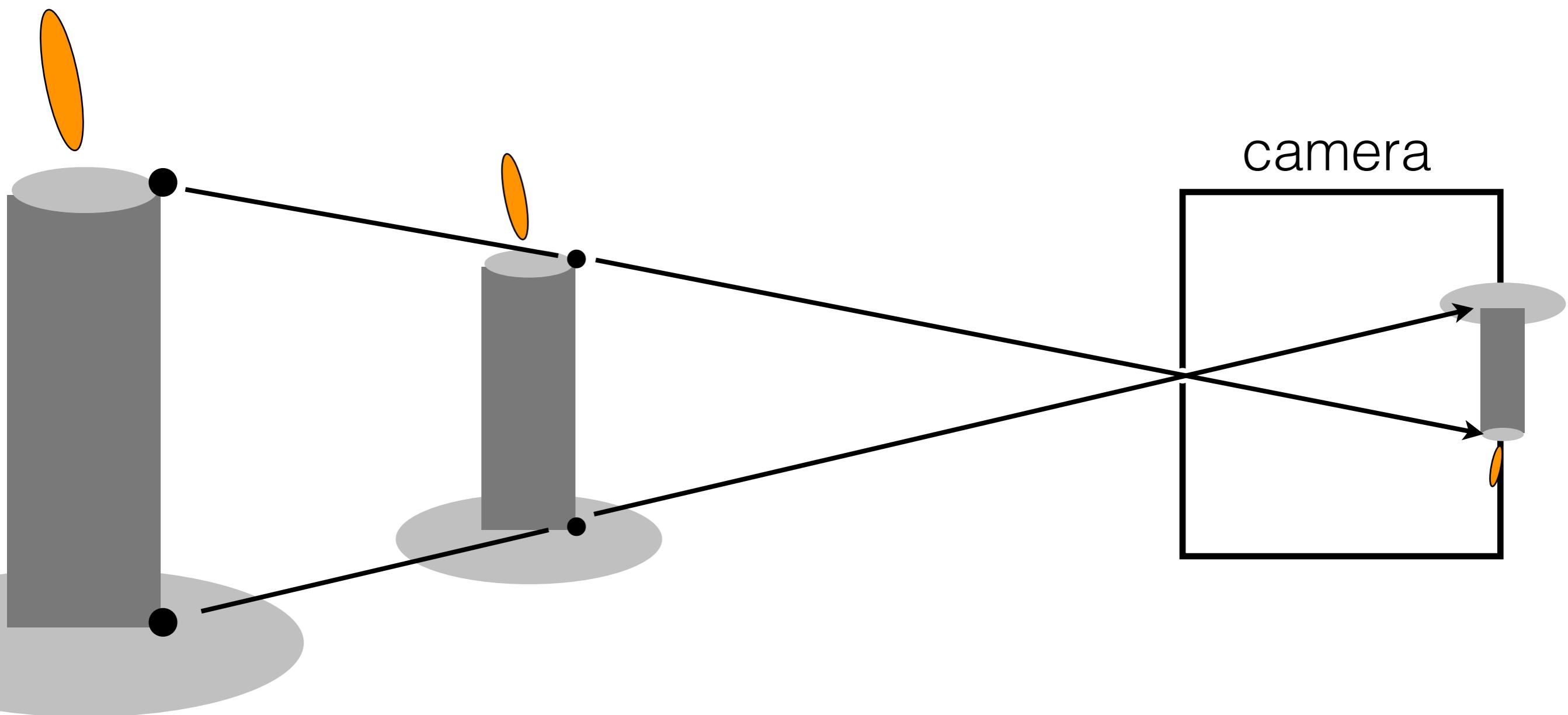
Bidirectional reflectance distribution function (BRDF)



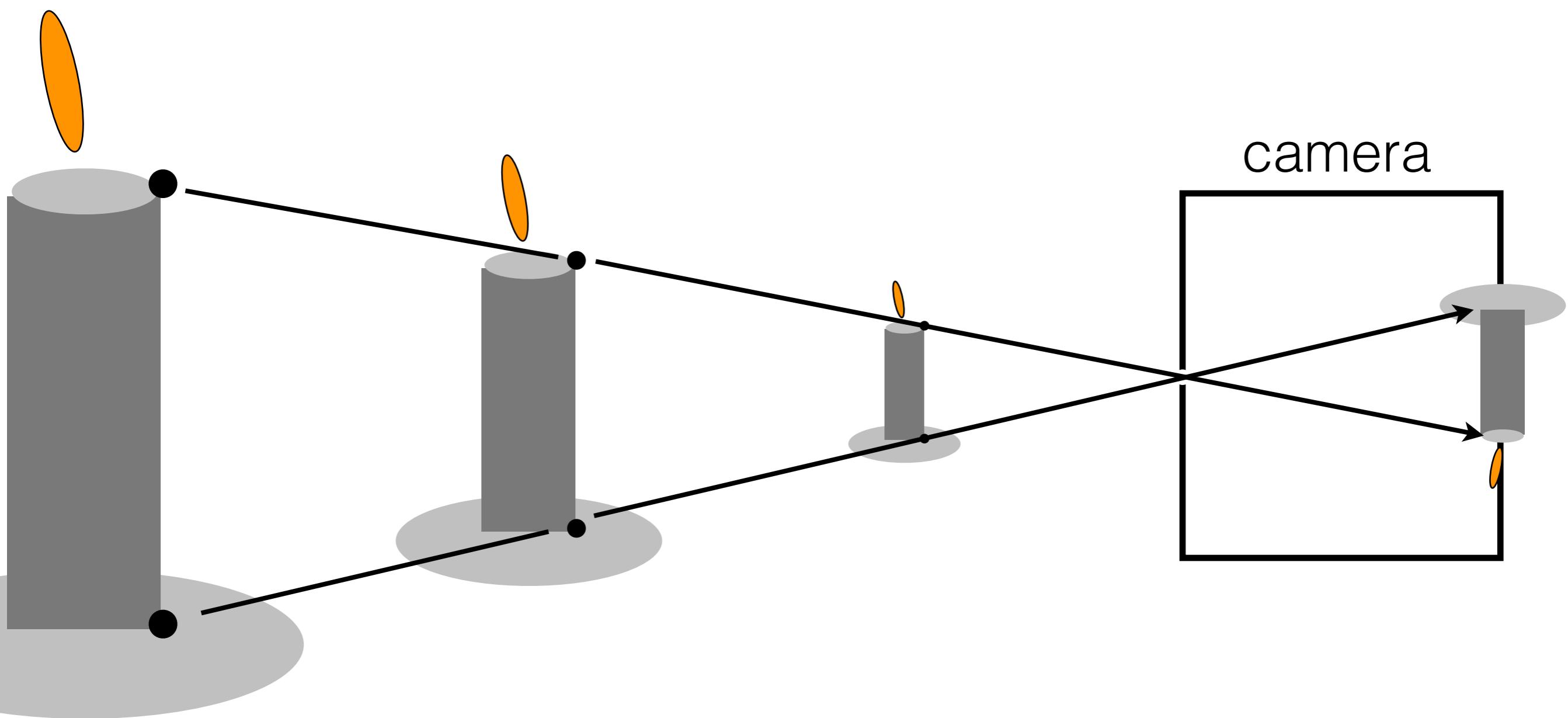
# Many worlds, one image



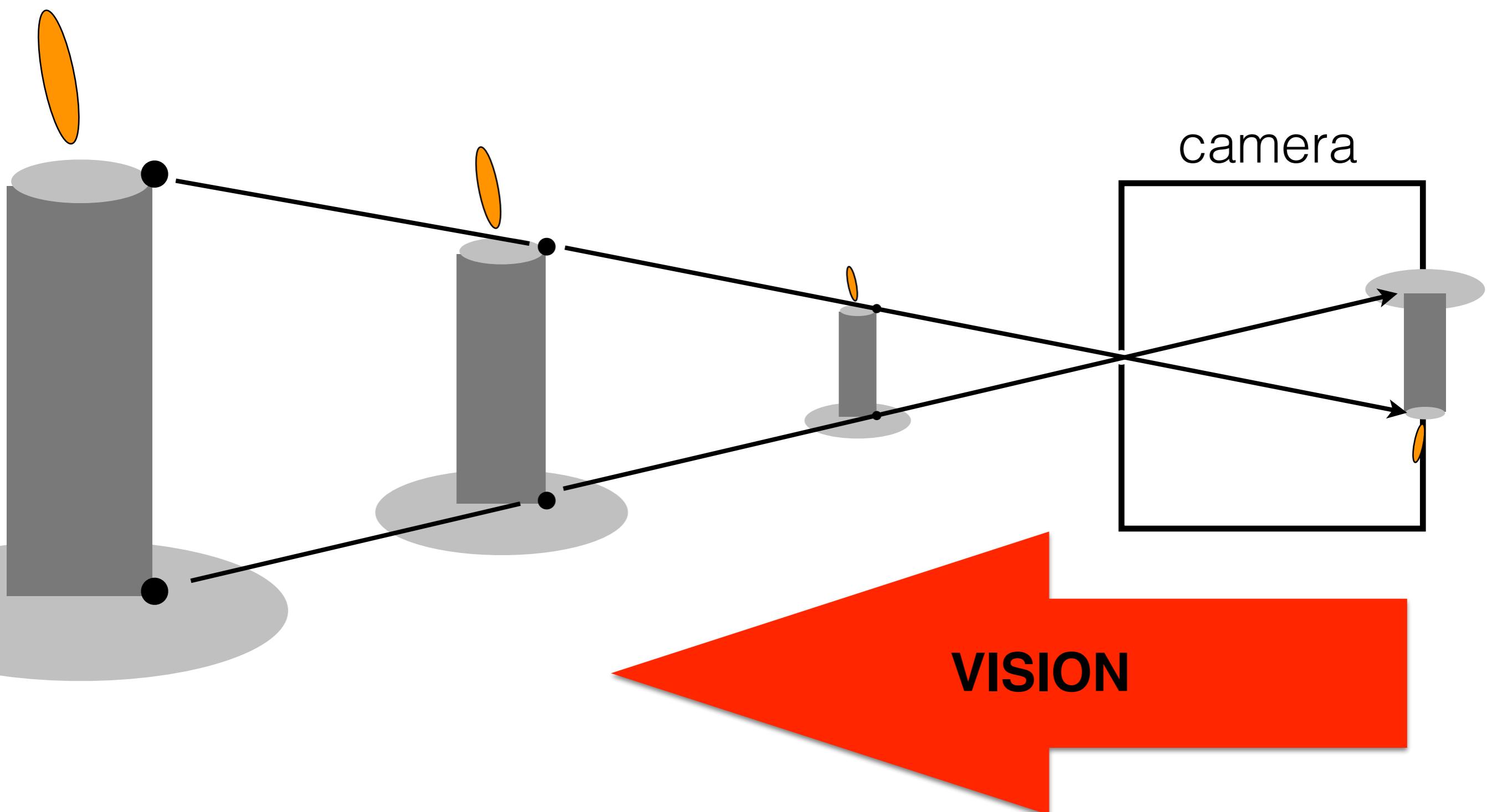
# Many worlds, one image



# Many worlds, one image



# Many worlds, one image







Shigeo Fukuda - Master of Deception  
Underground Piano

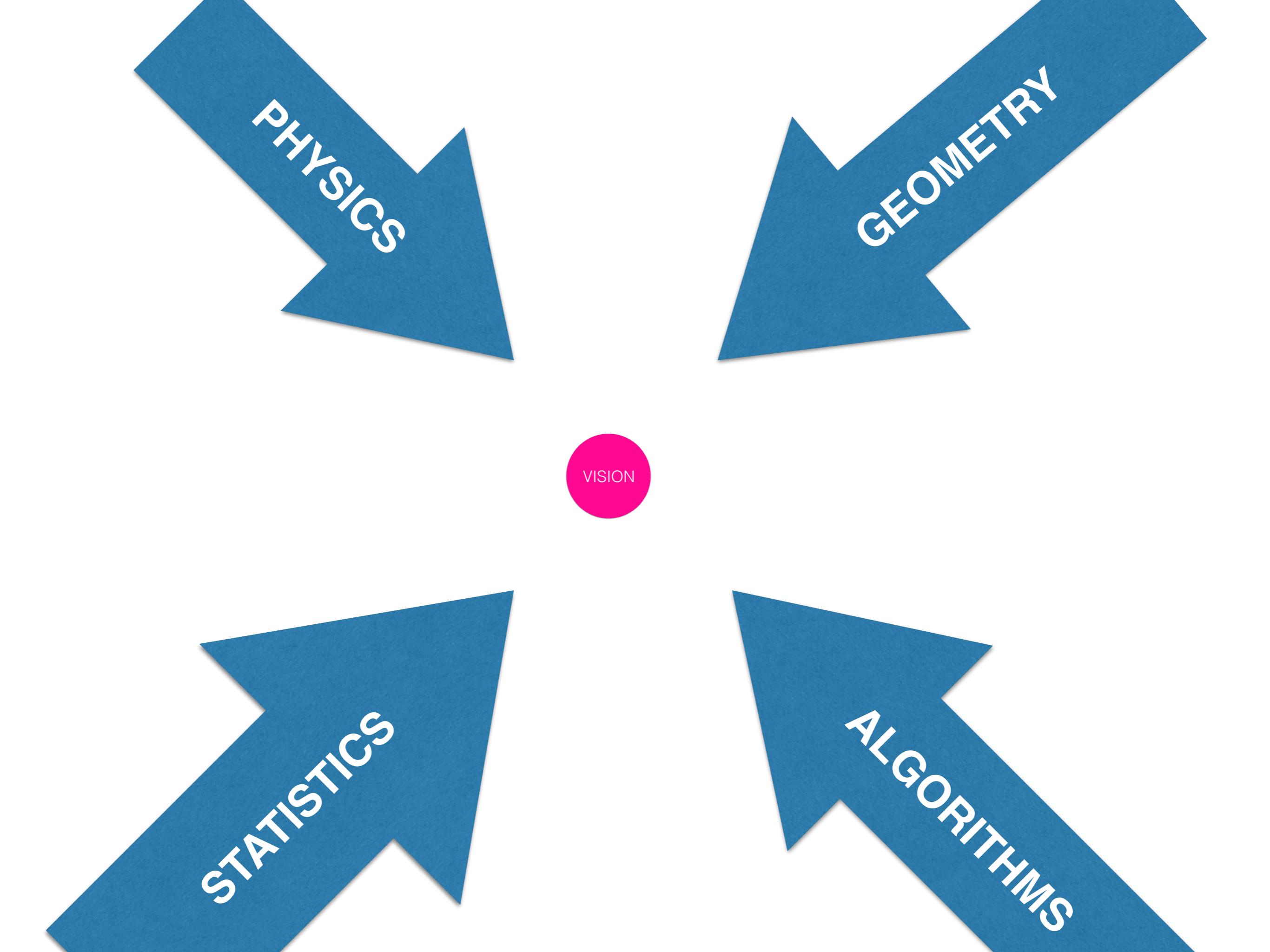
VISION

PHYSICS

GEOMETRY

STATISTICS

ALGORITHMS



PHYSICS

GEOMETRY

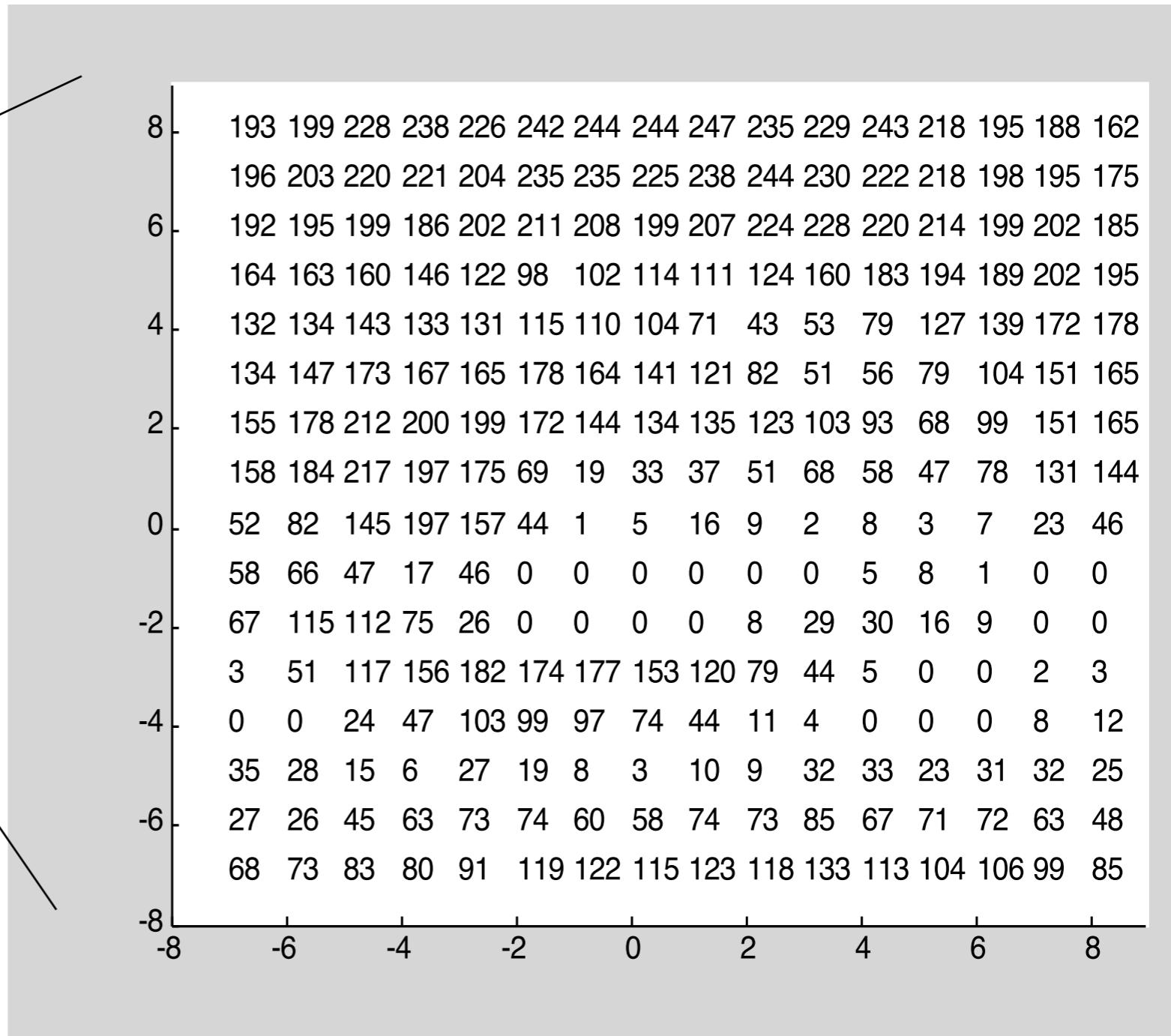
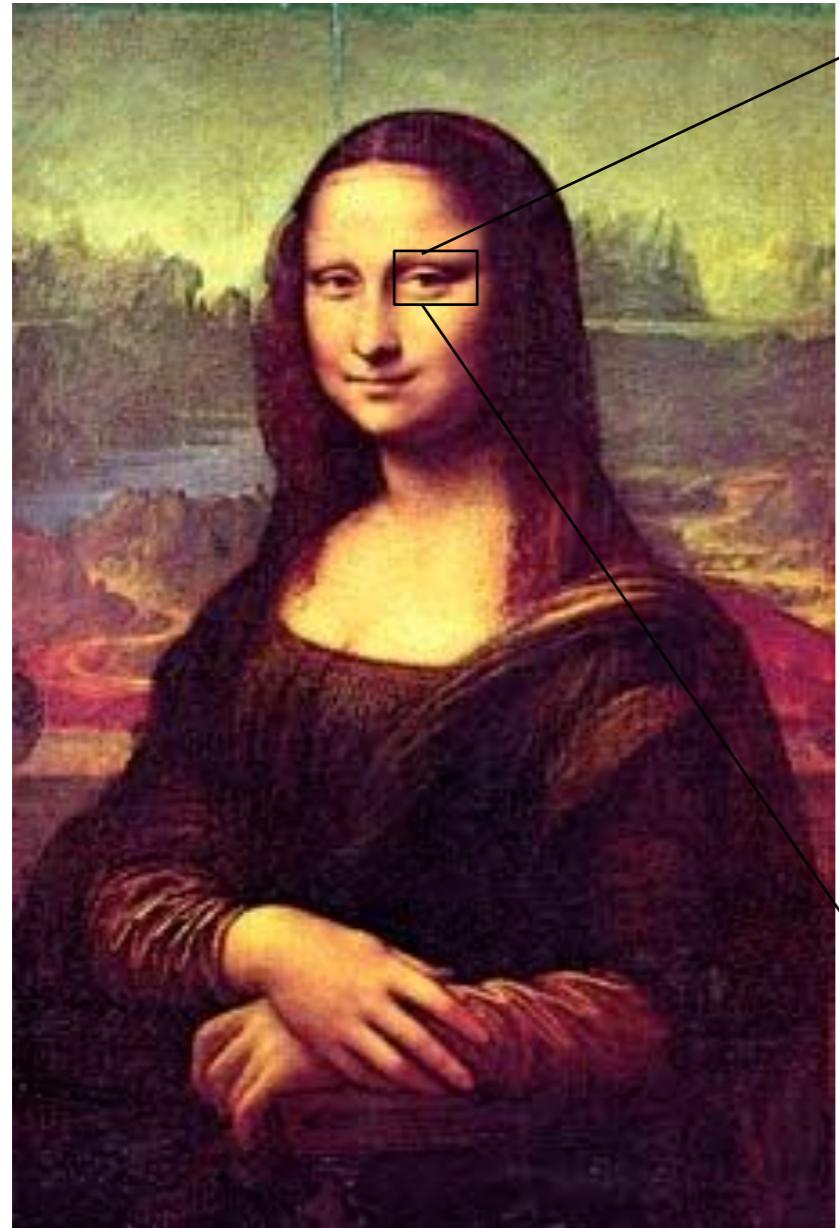
STATISTICS

ALGORITHMS

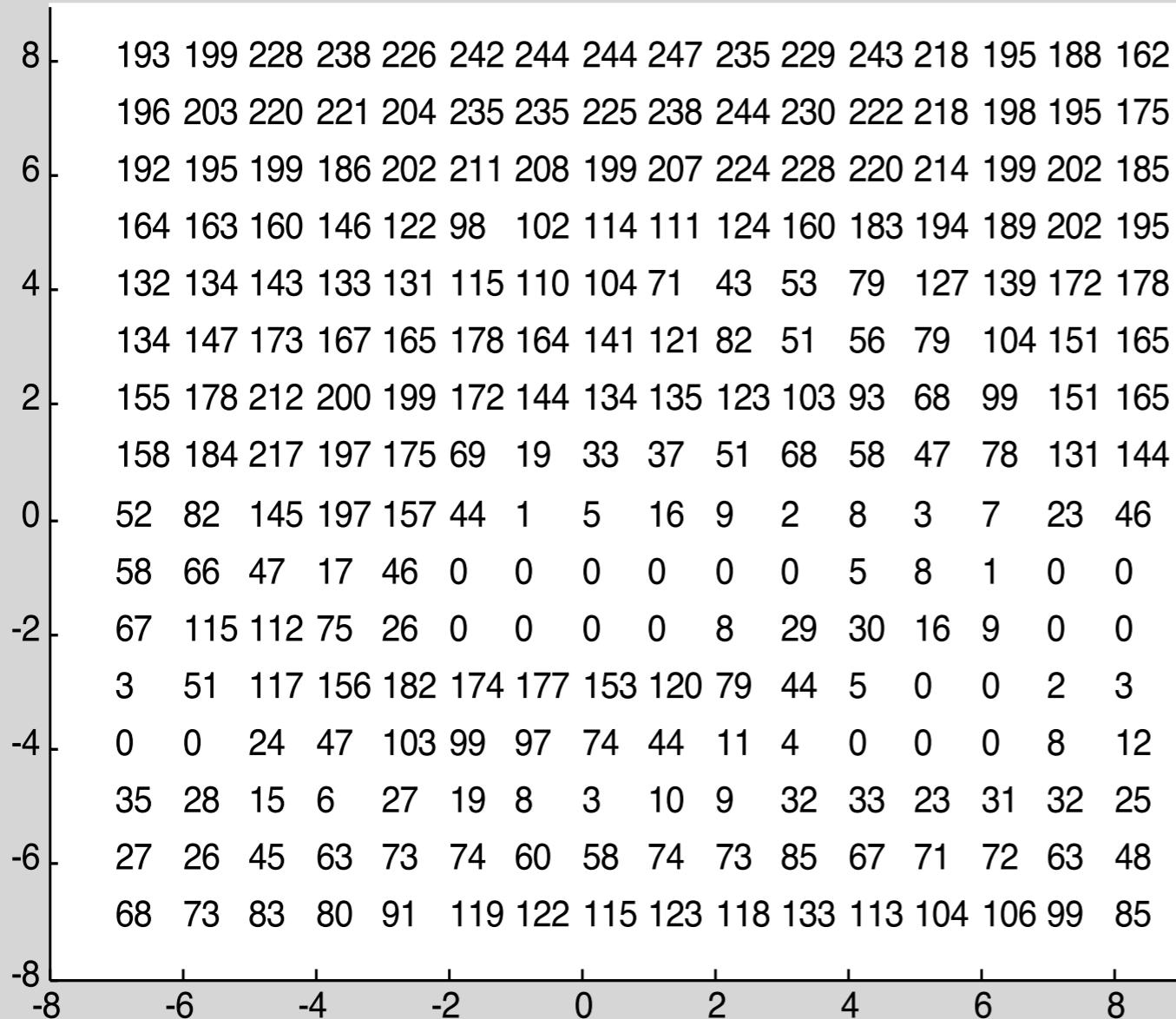
VISION

A more modest proposal...

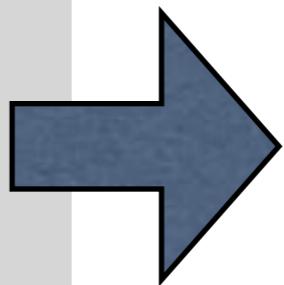
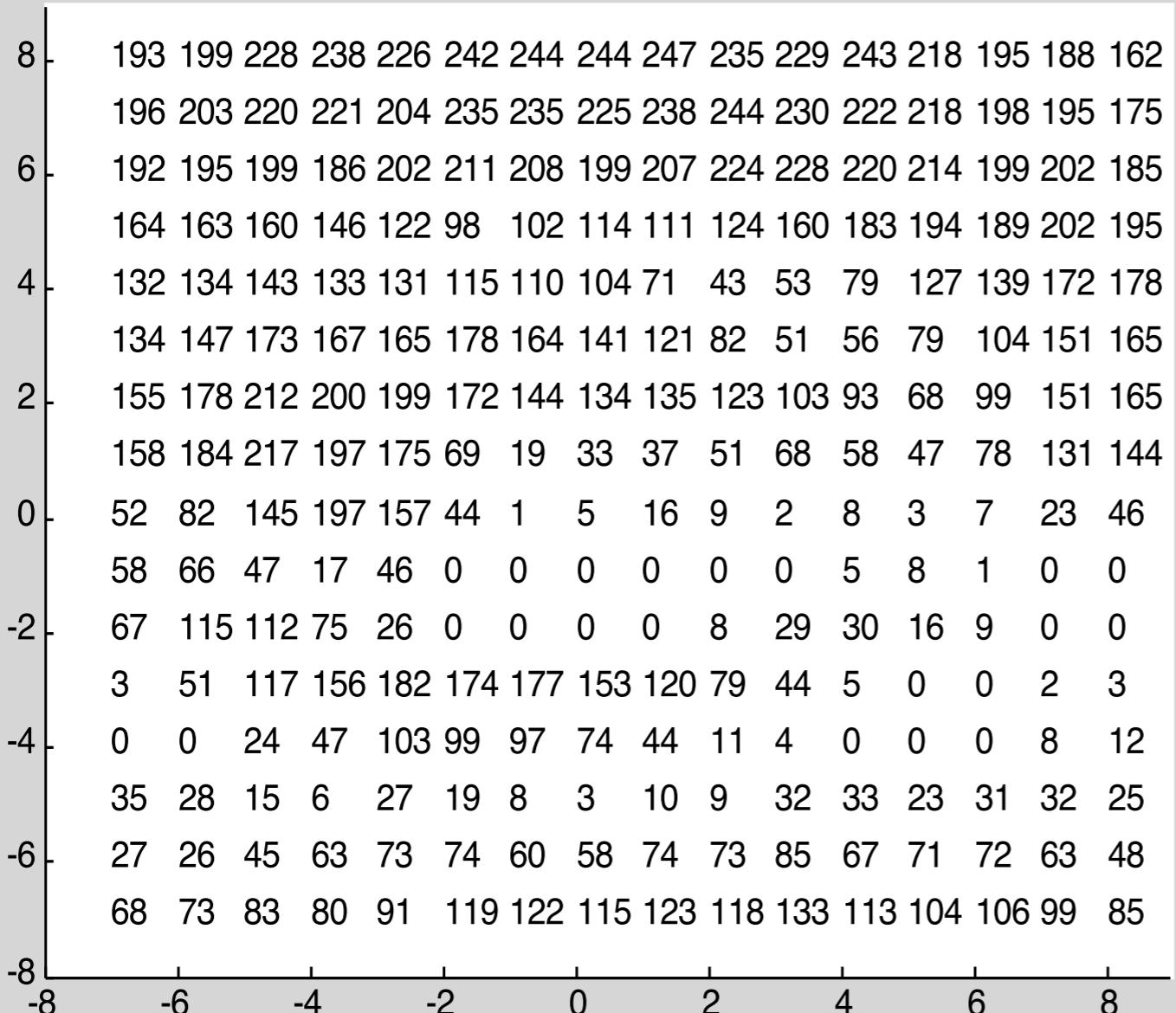
# What is in an image?



# The challenge

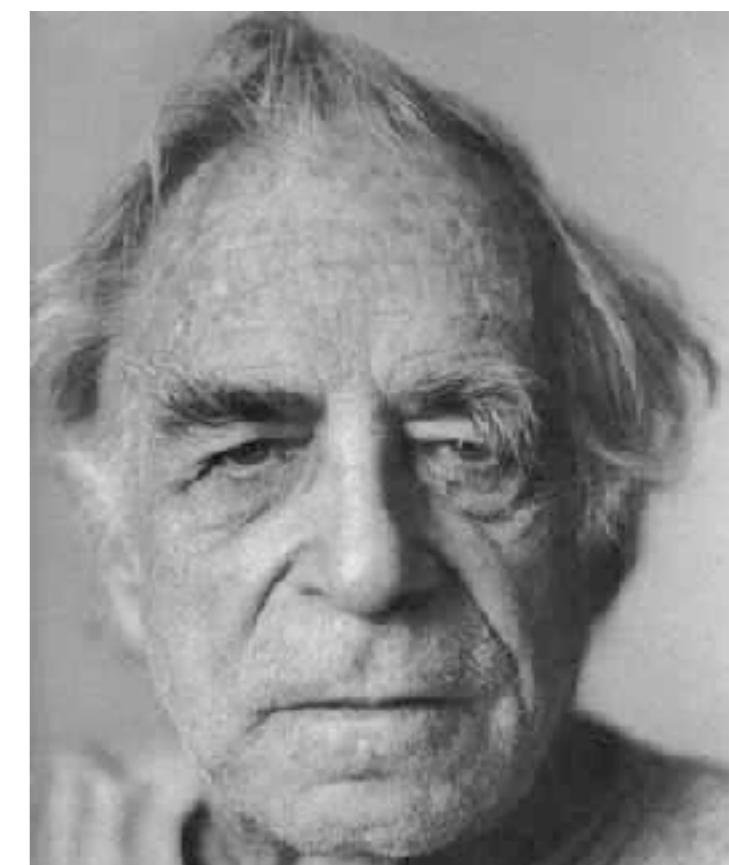
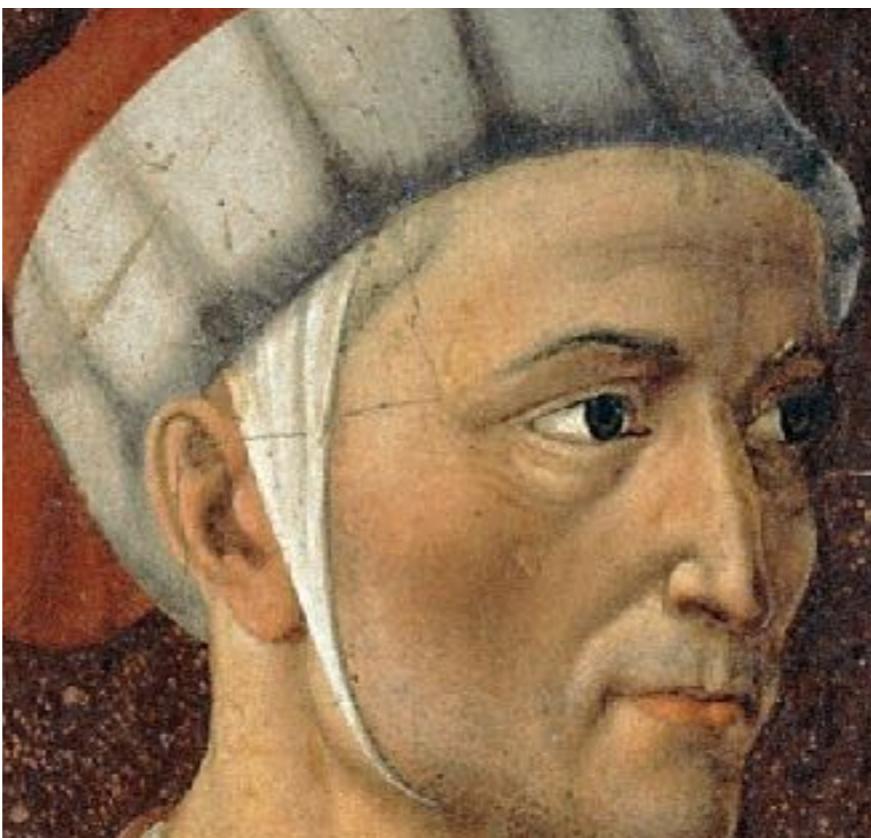
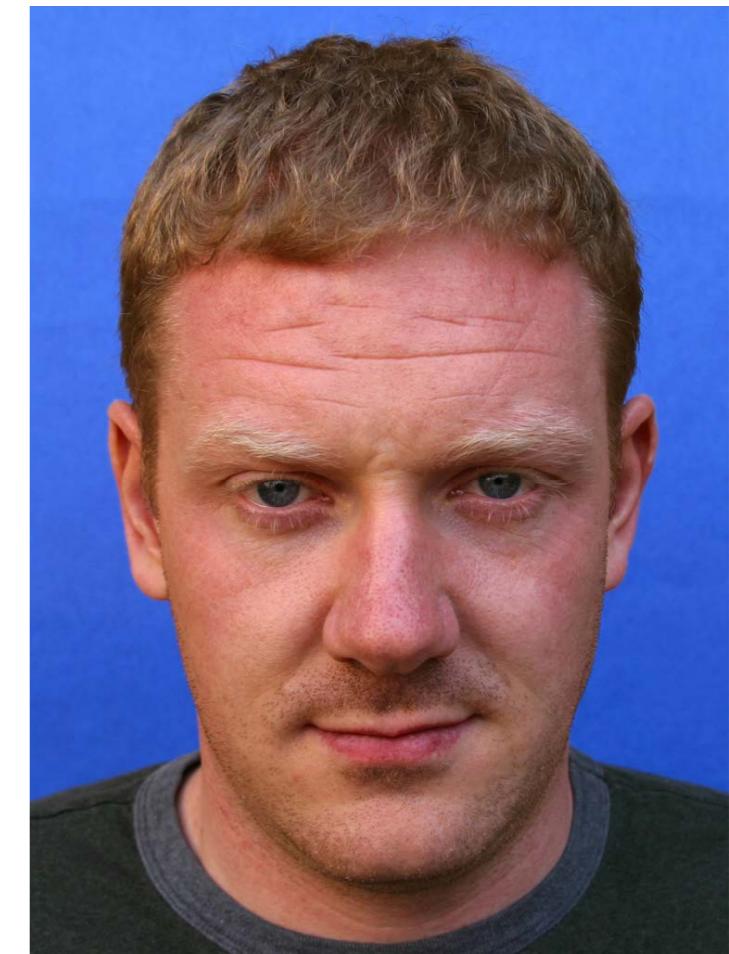
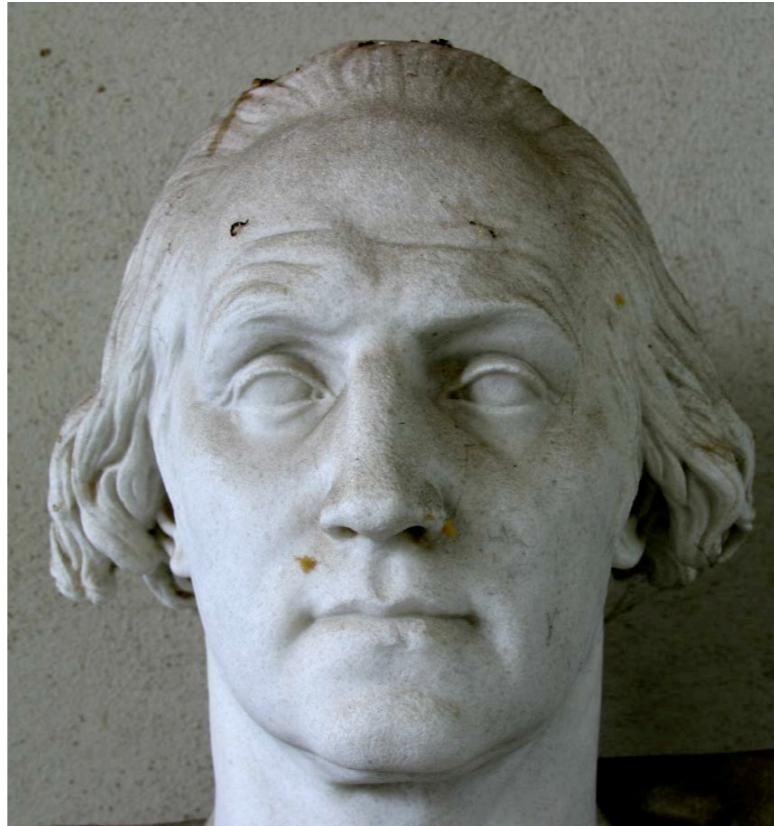


# The challenge

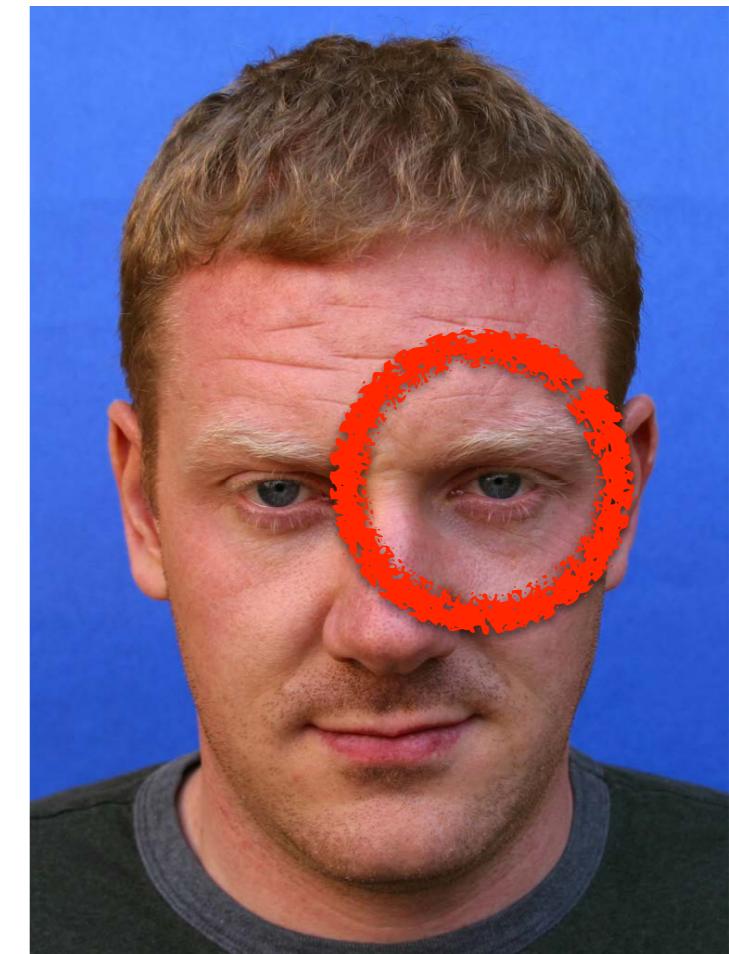
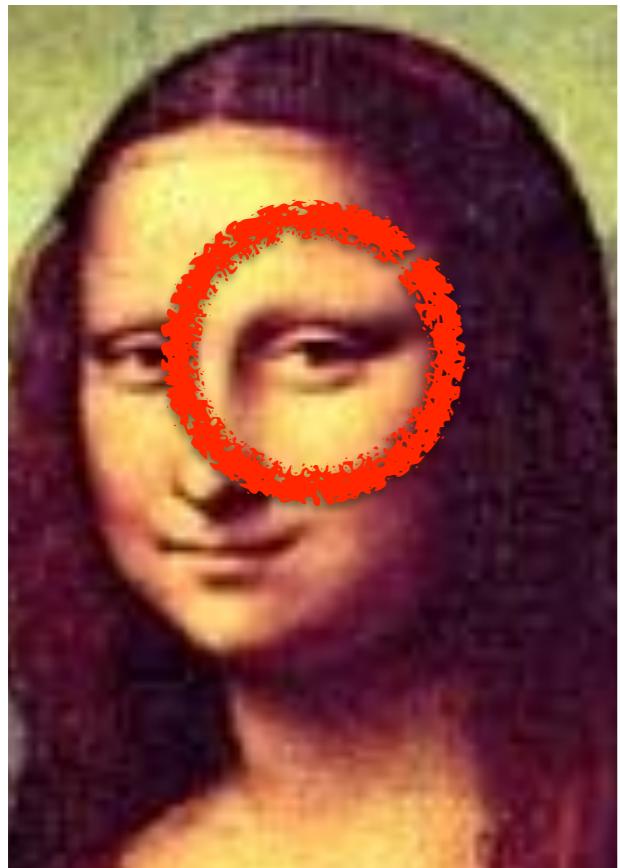


“Eye”

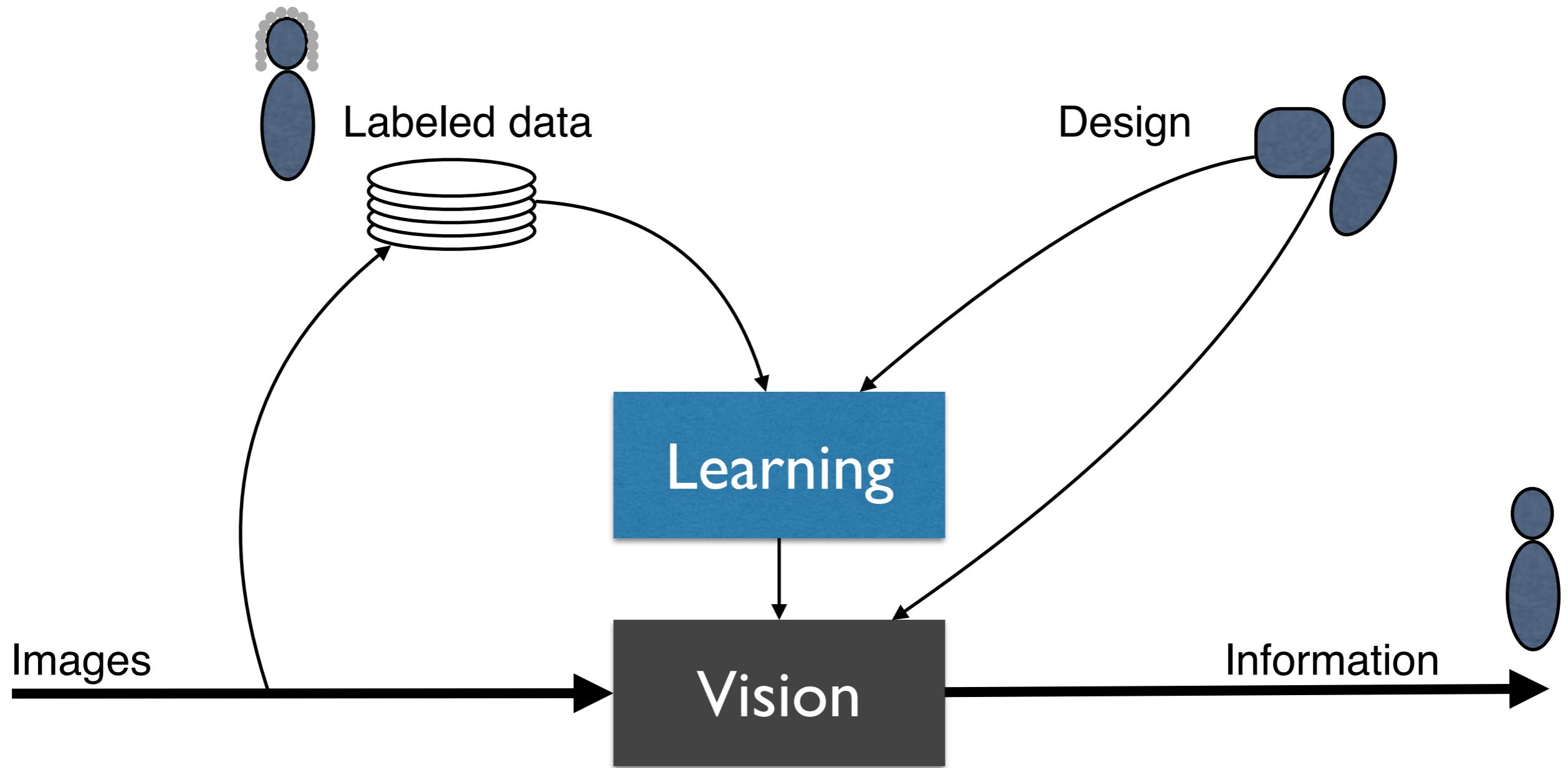
# “Eye”

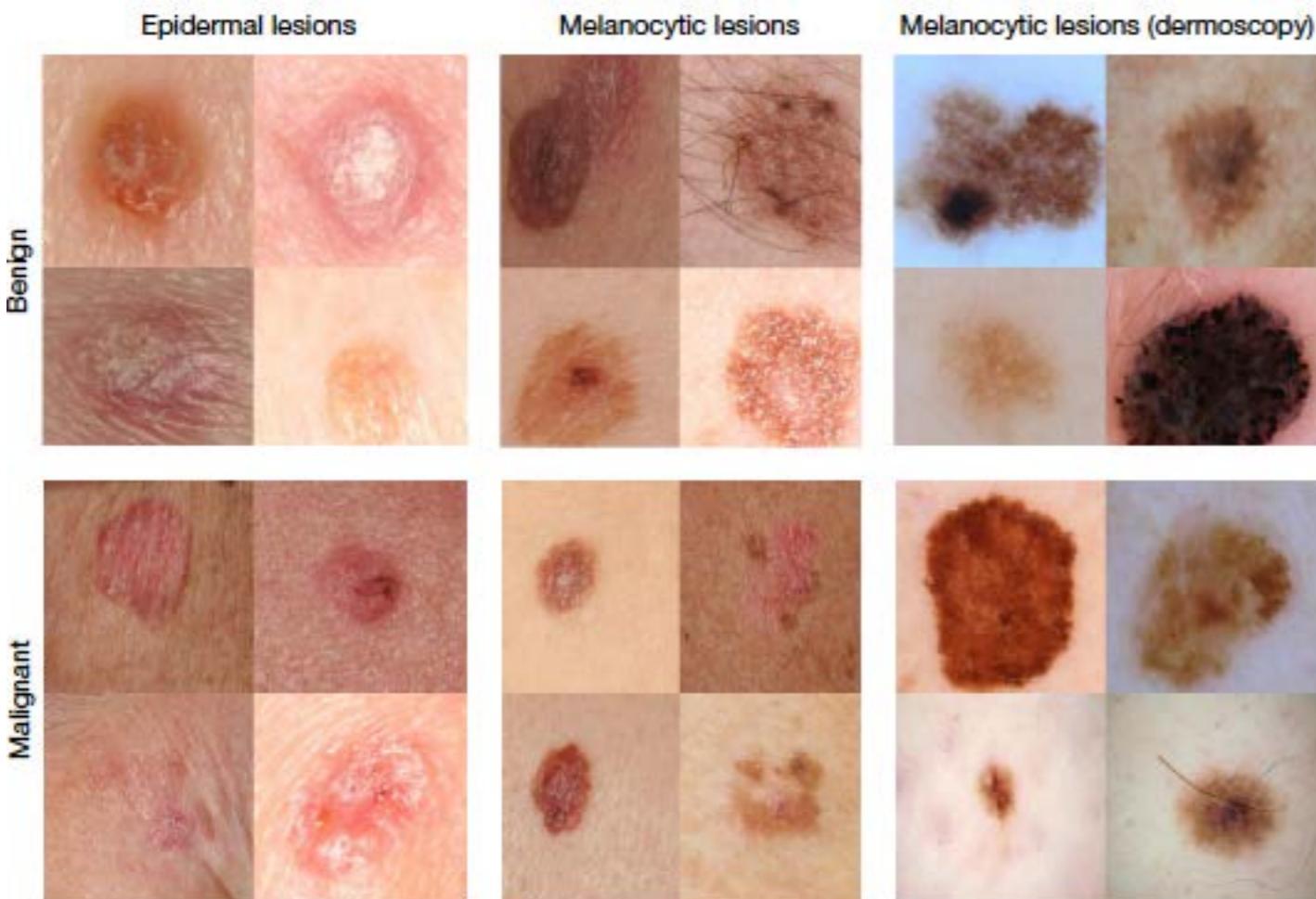


# “Eye”

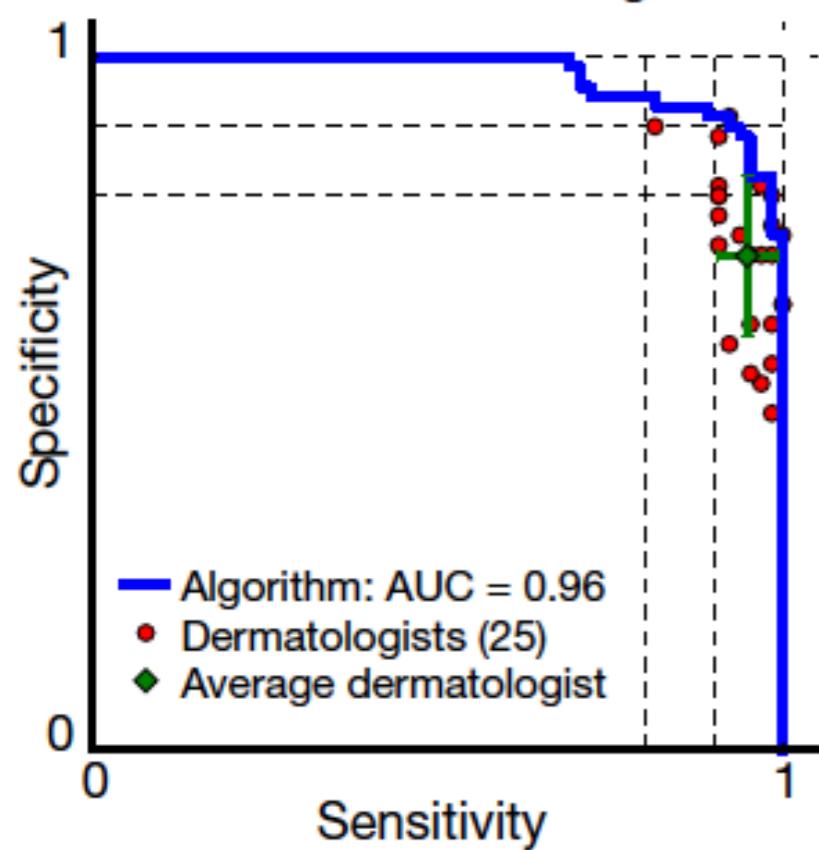


# Learning-based approach

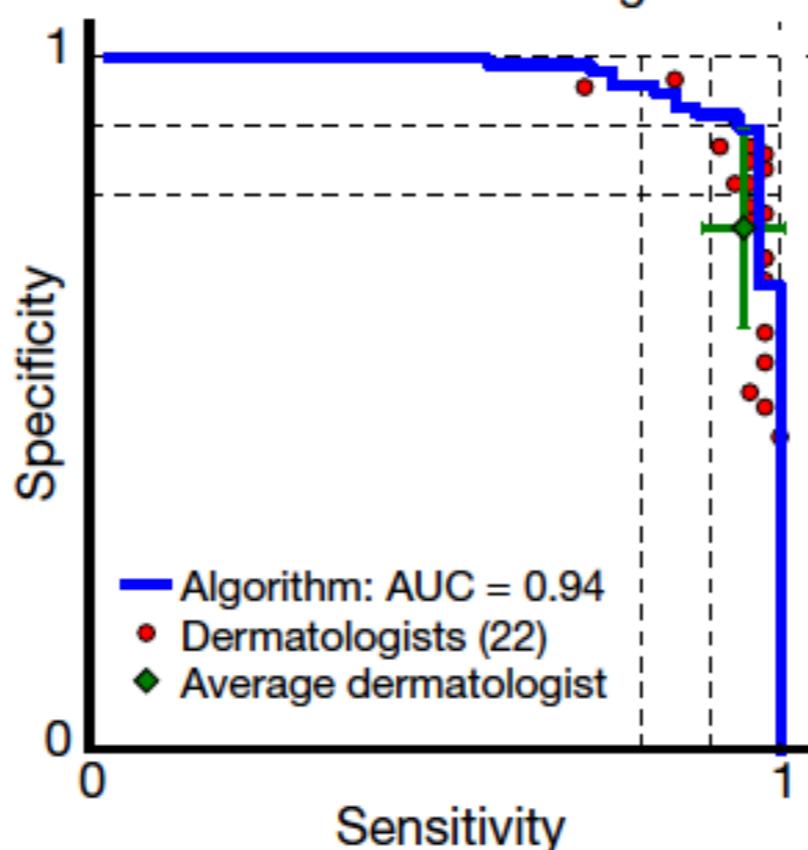




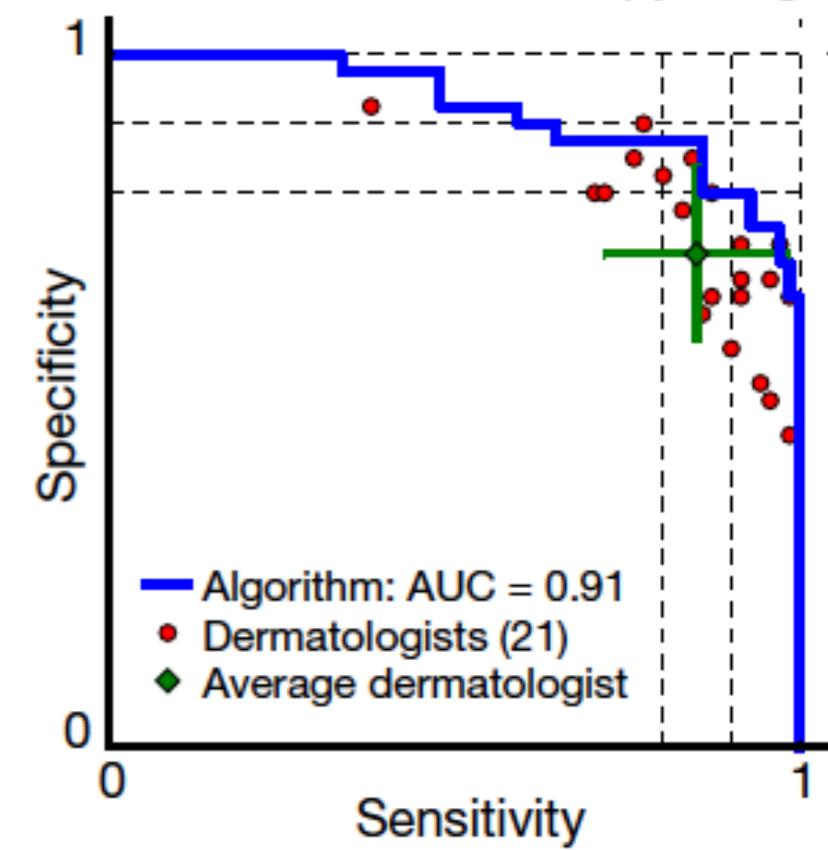
Carcinoma: 135 images



Melanoma: 130 images



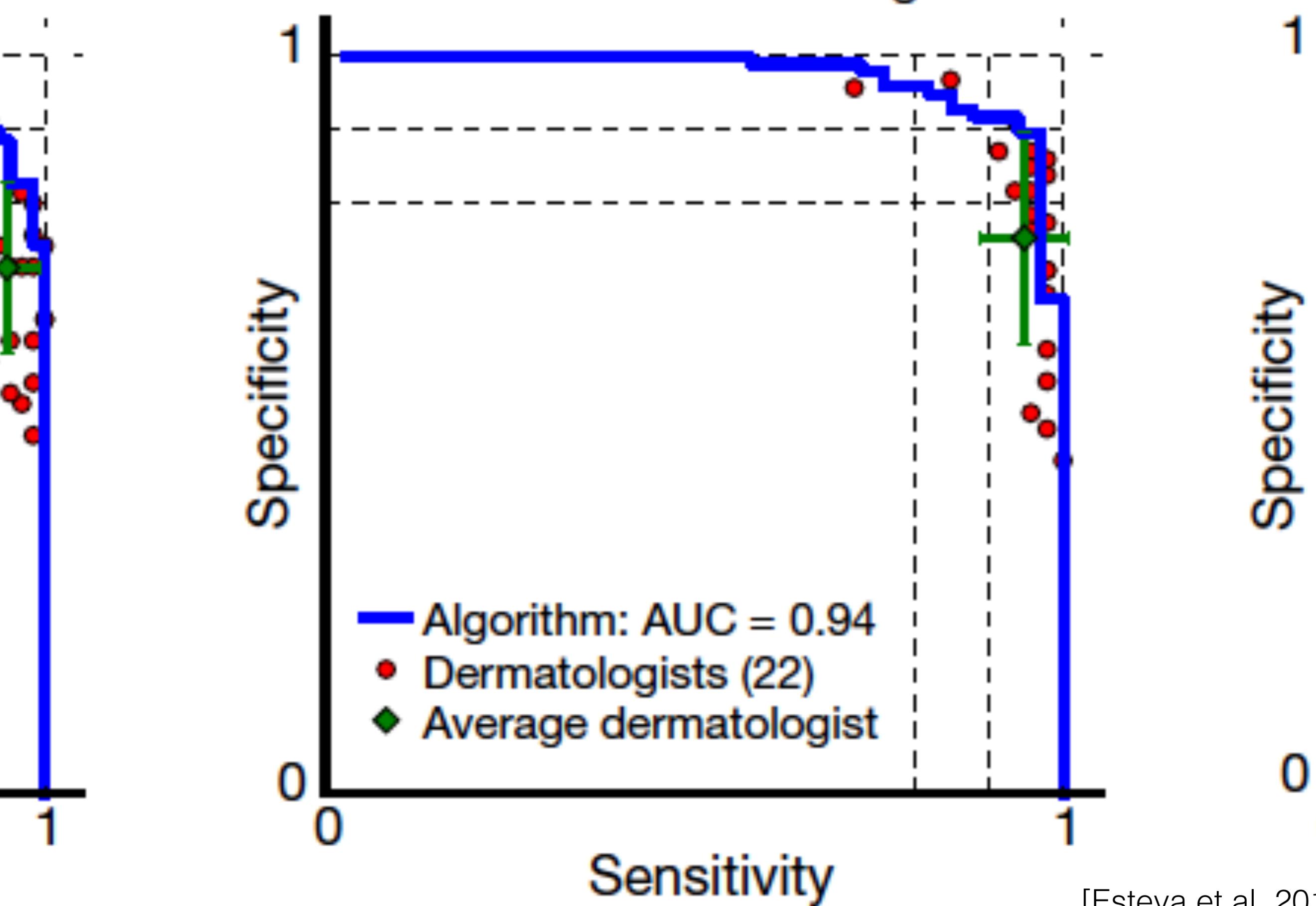
Melanoma: 111 dermoscopy images

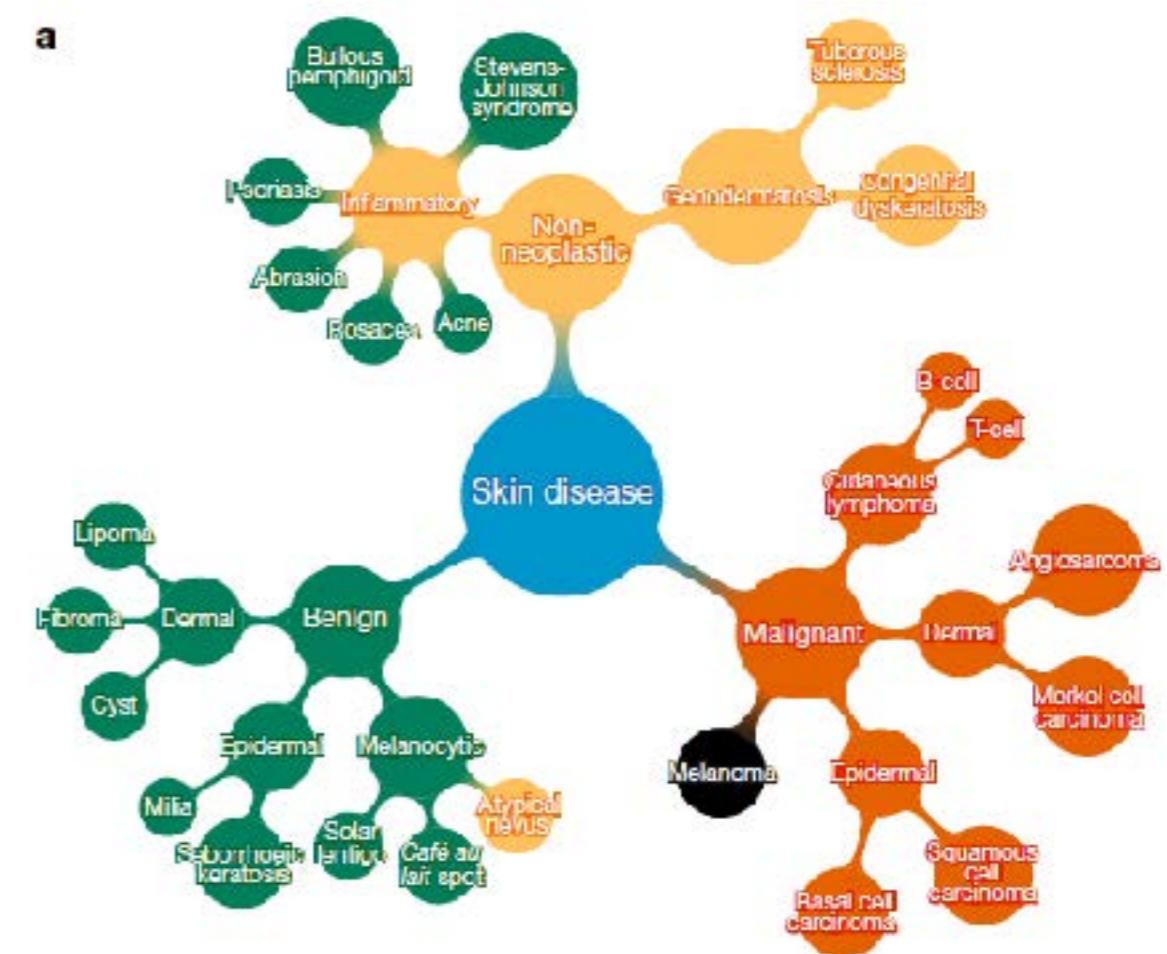
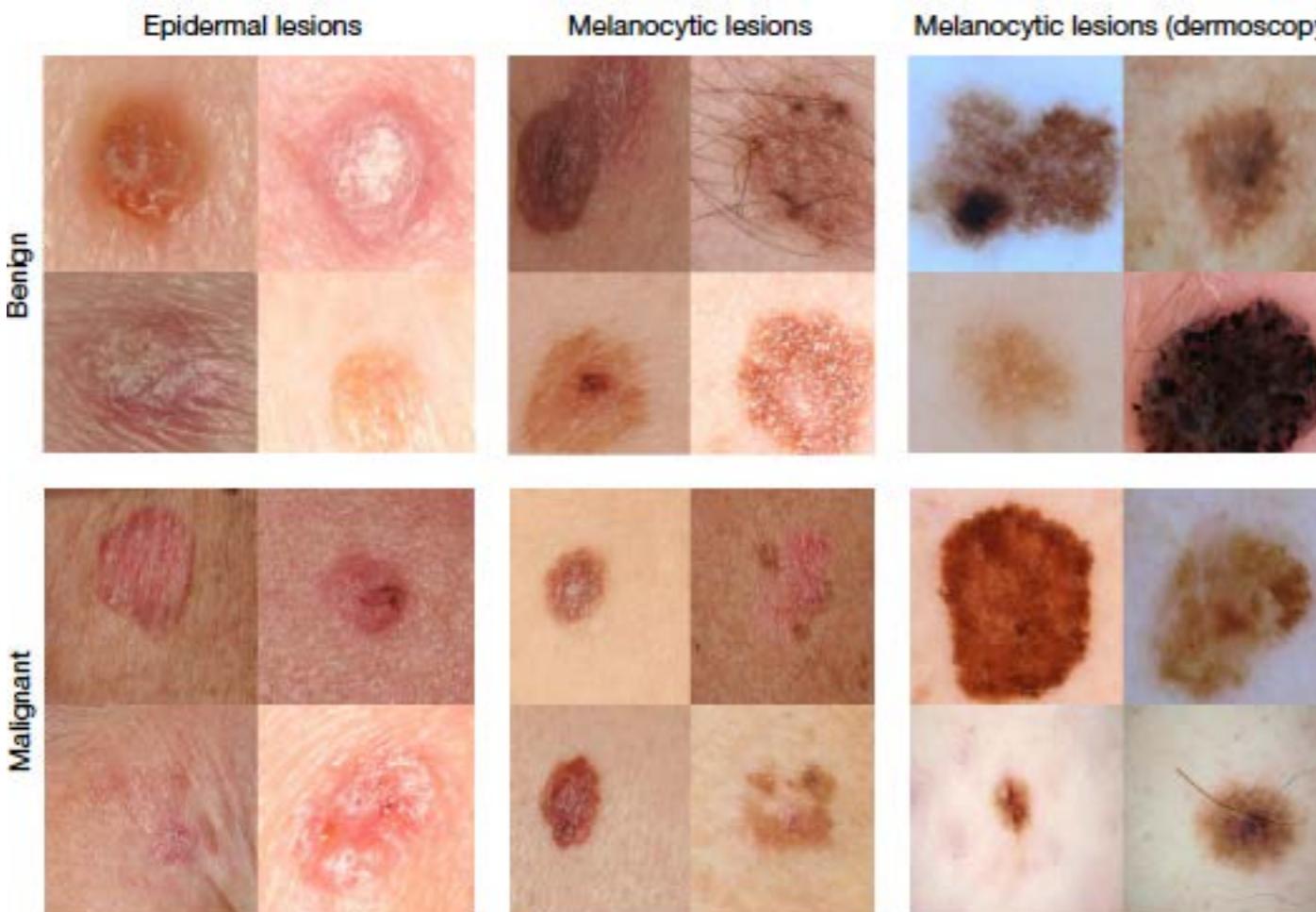


[Esteva et al. 2017]

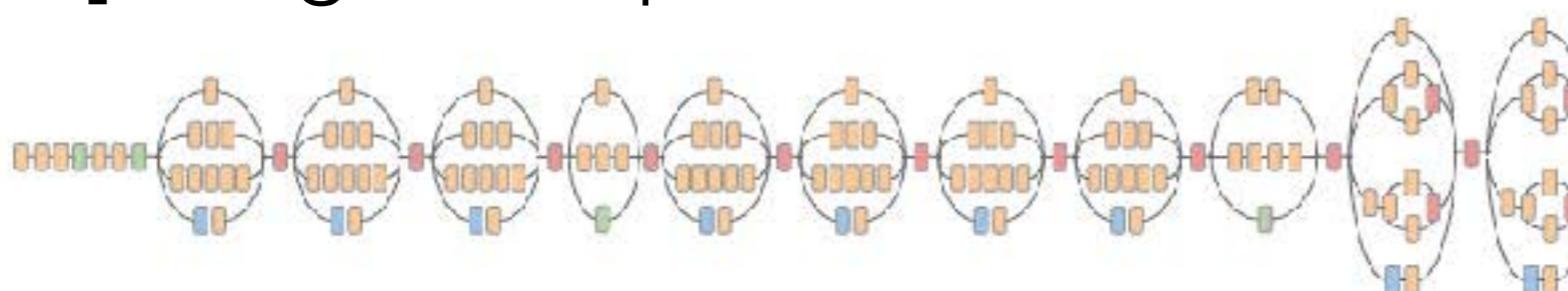
Mela

## Melanoma: 130 images

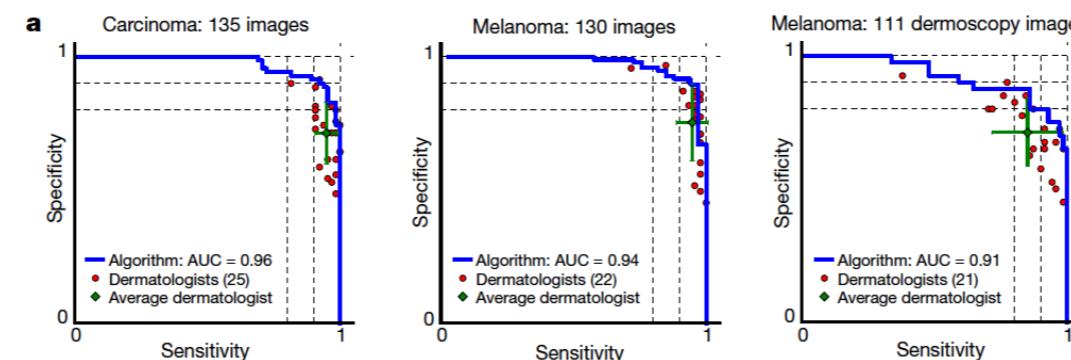




## [Google inception architecture 2015]

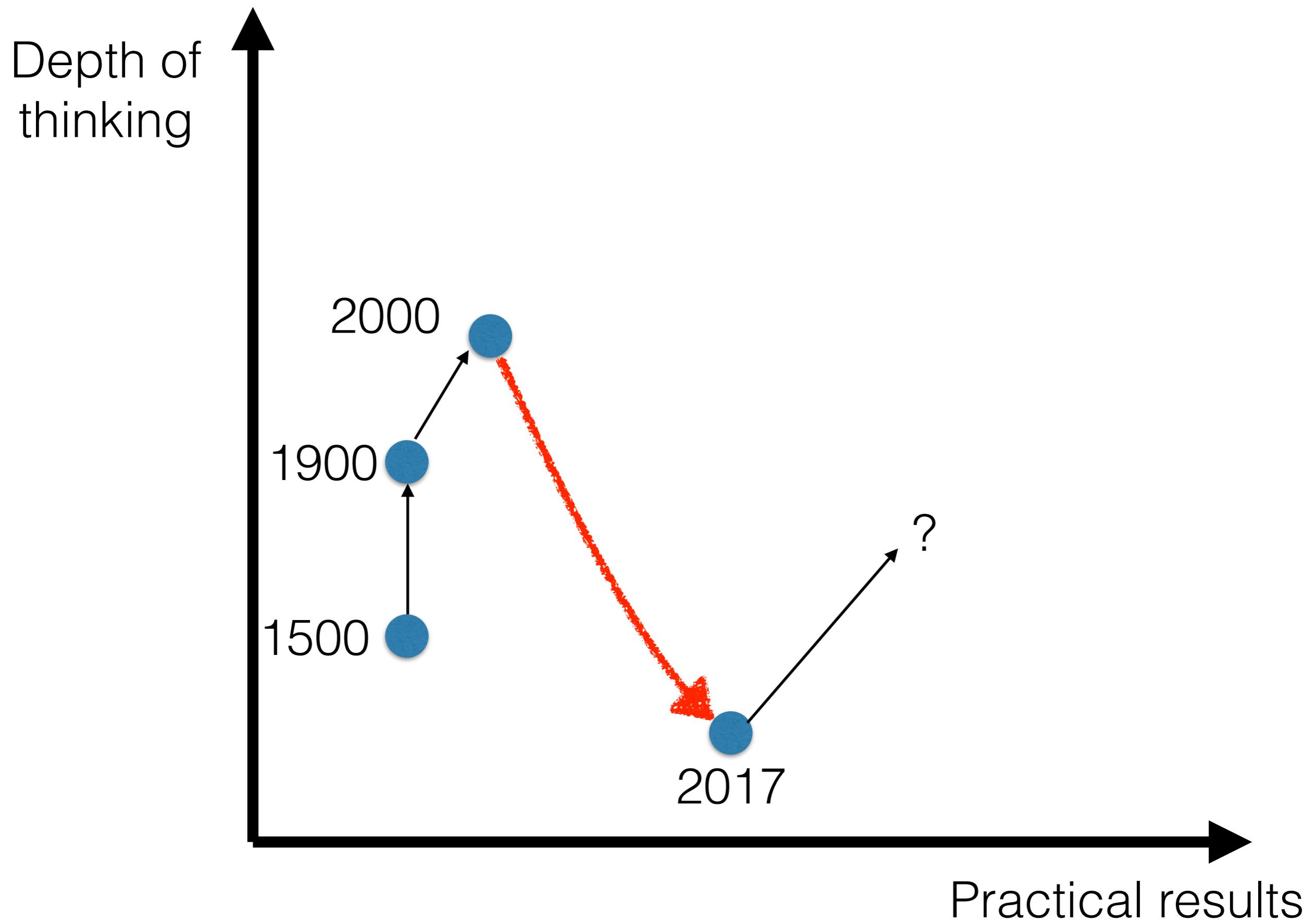


- Convolution
- AvgPool
- MaxPool
- Concat
- Dropout
- Fully connected
- Softmax



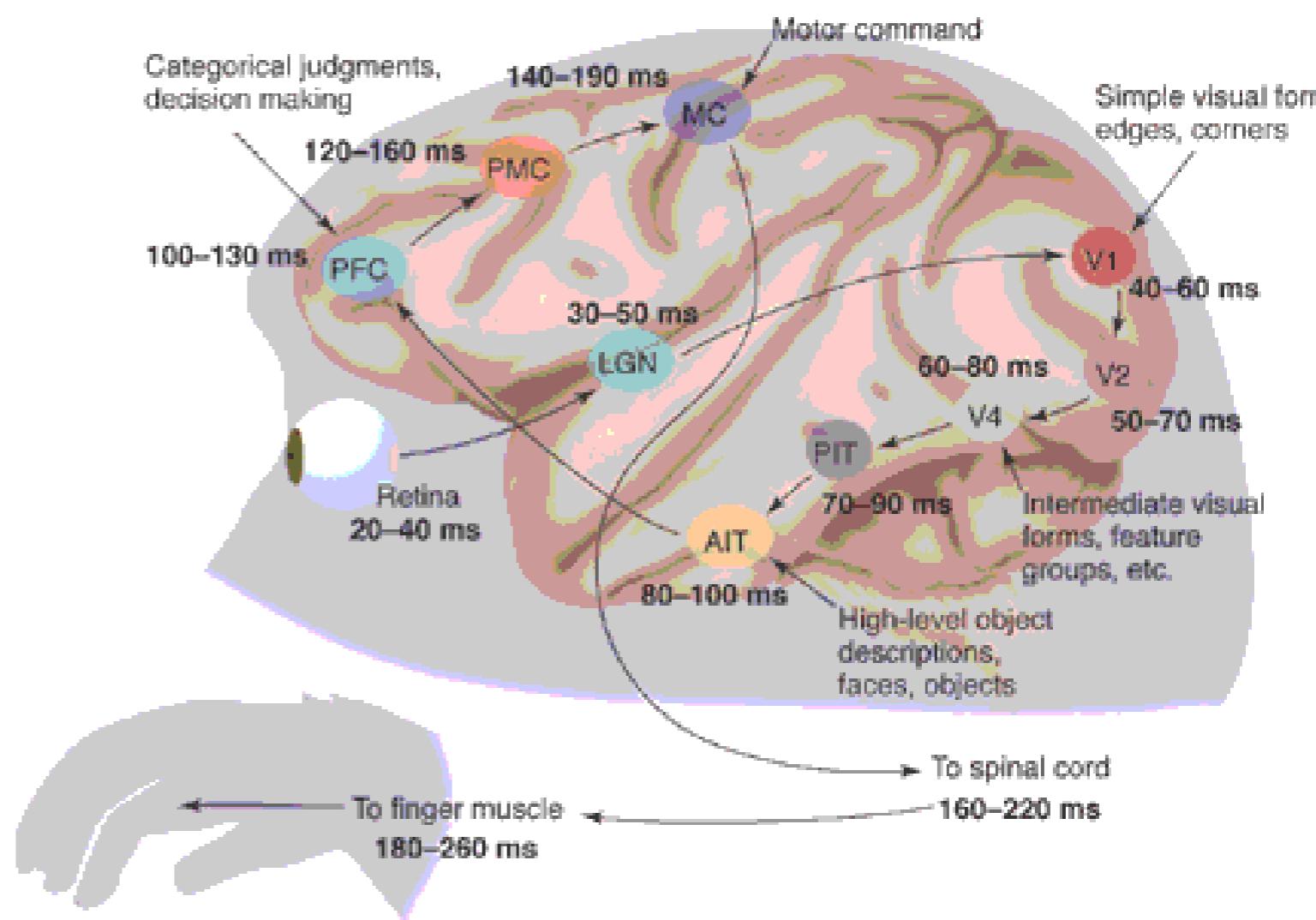
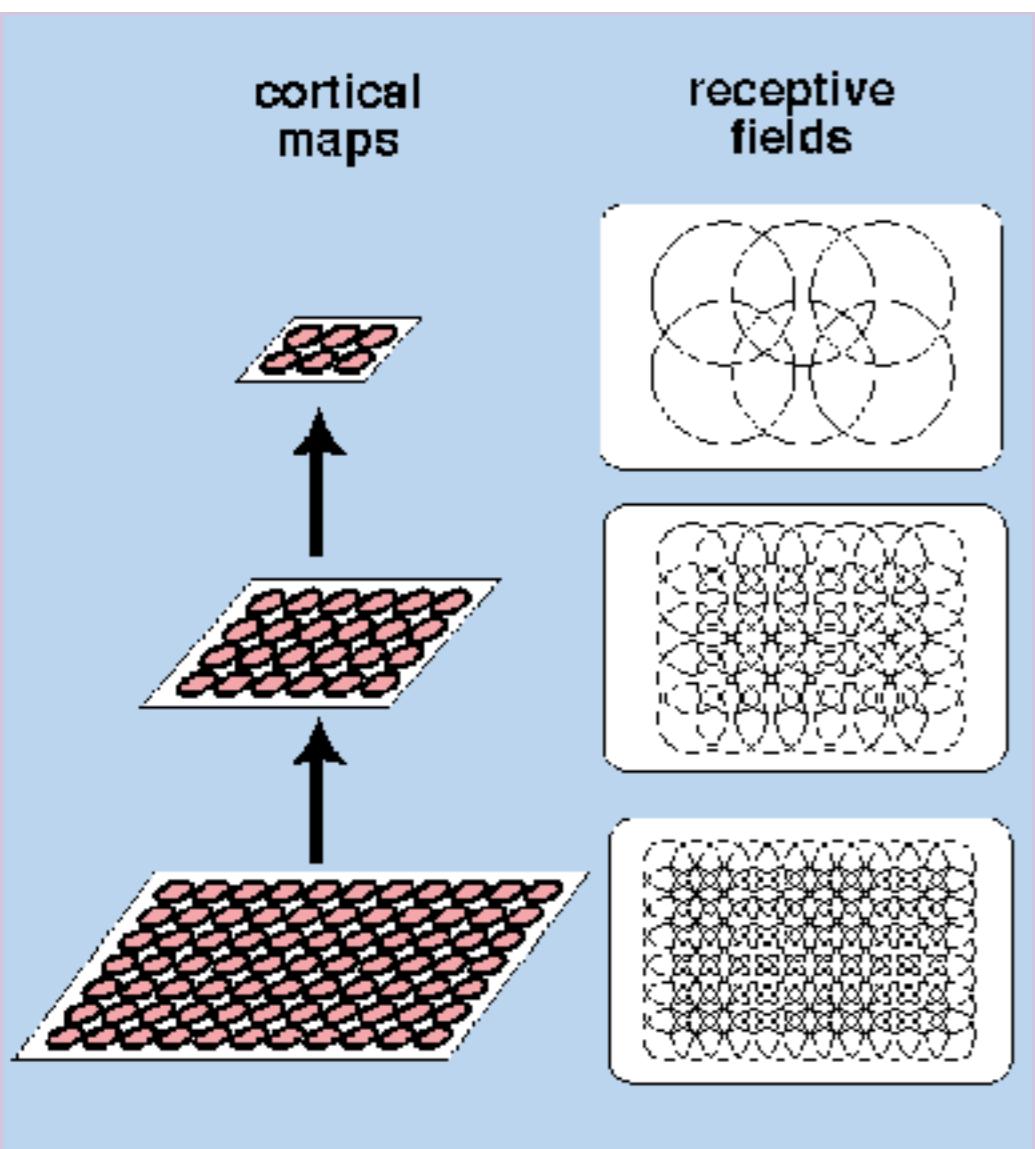
- Acral-lentiginous melanoma
- Amelanotic melanoma
- Lentigo melanoma
- ...
- Blue nevus
- Halo nevus
- Mongolian spot
- ...

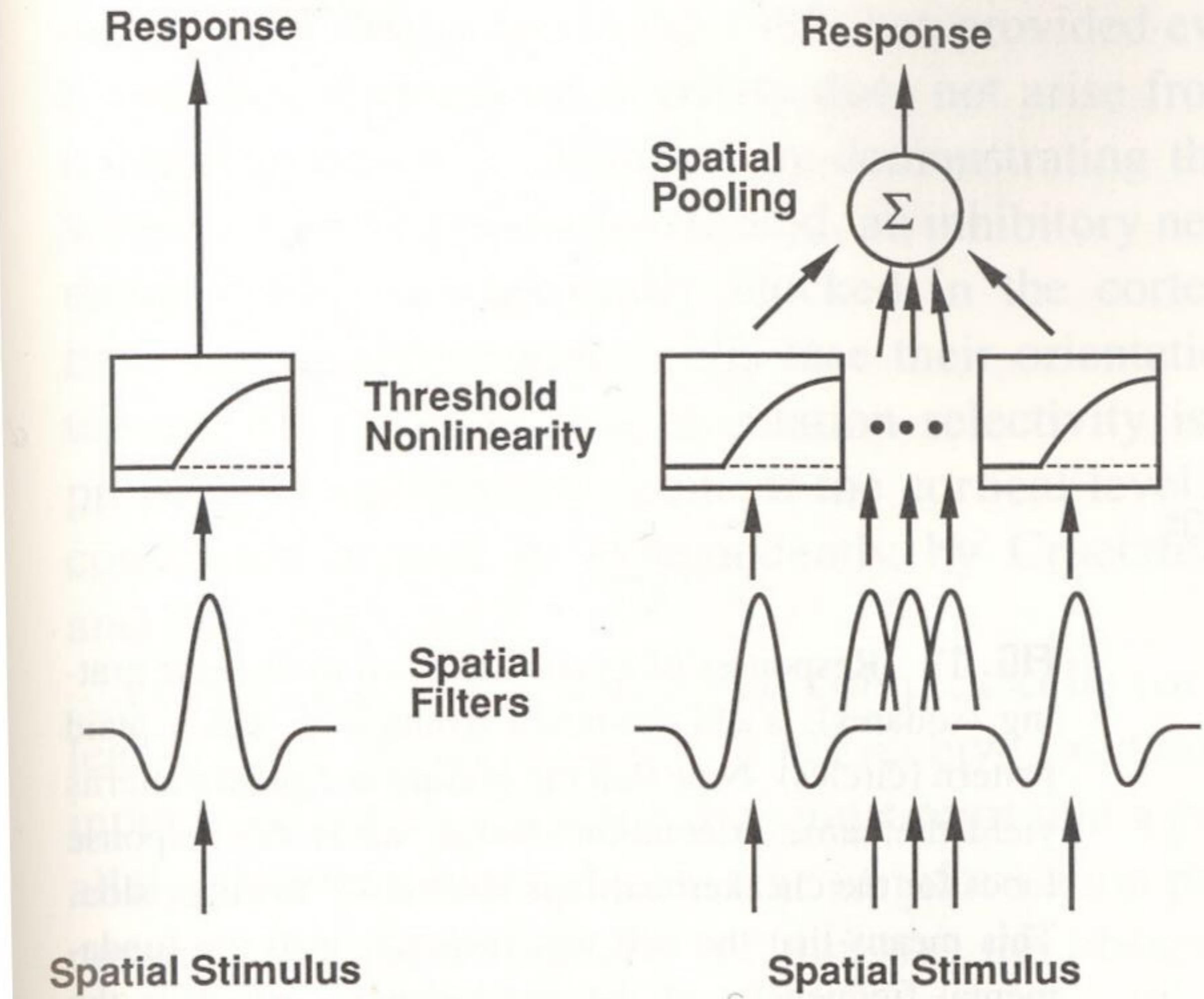
[Esteva et al. 2017]

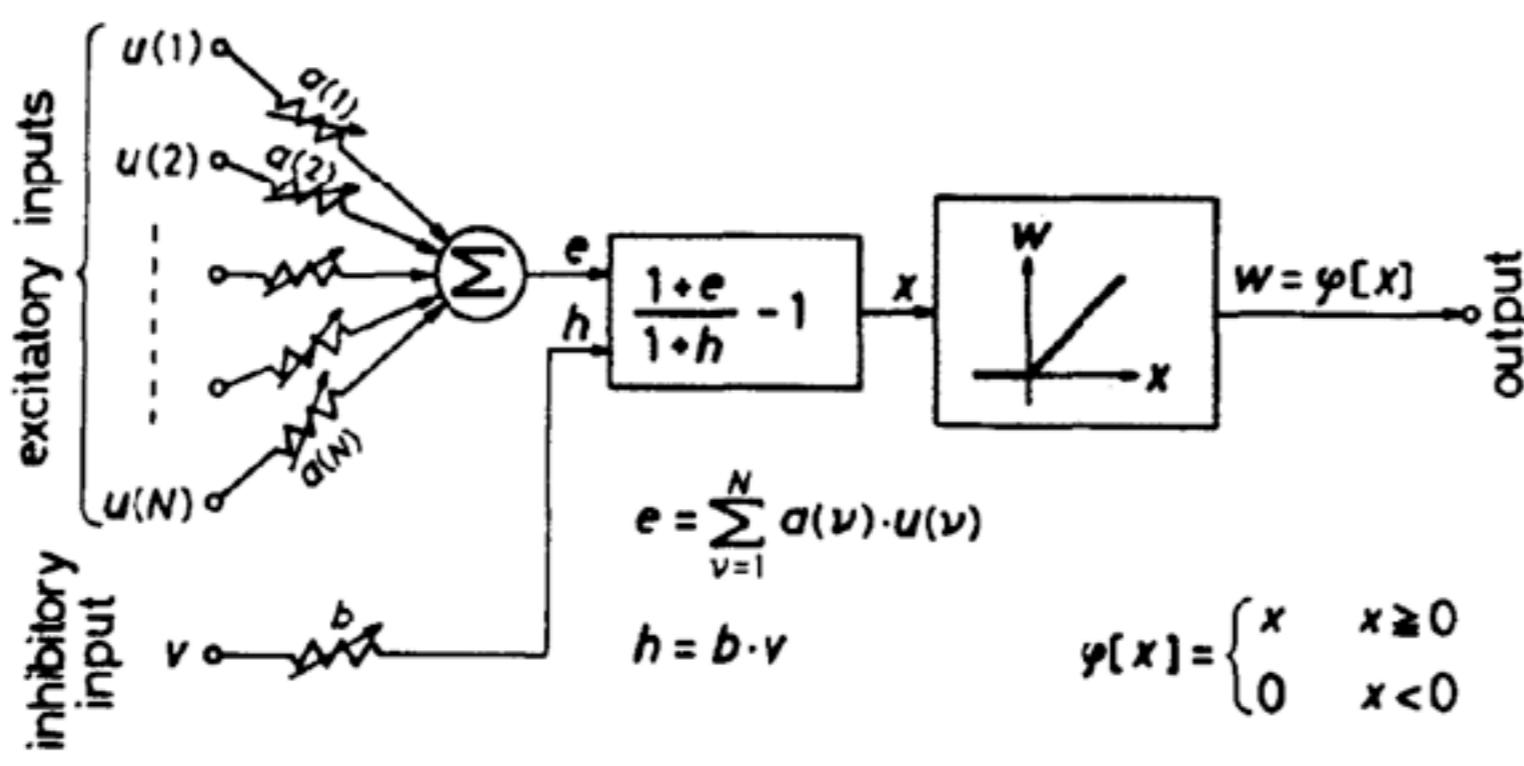
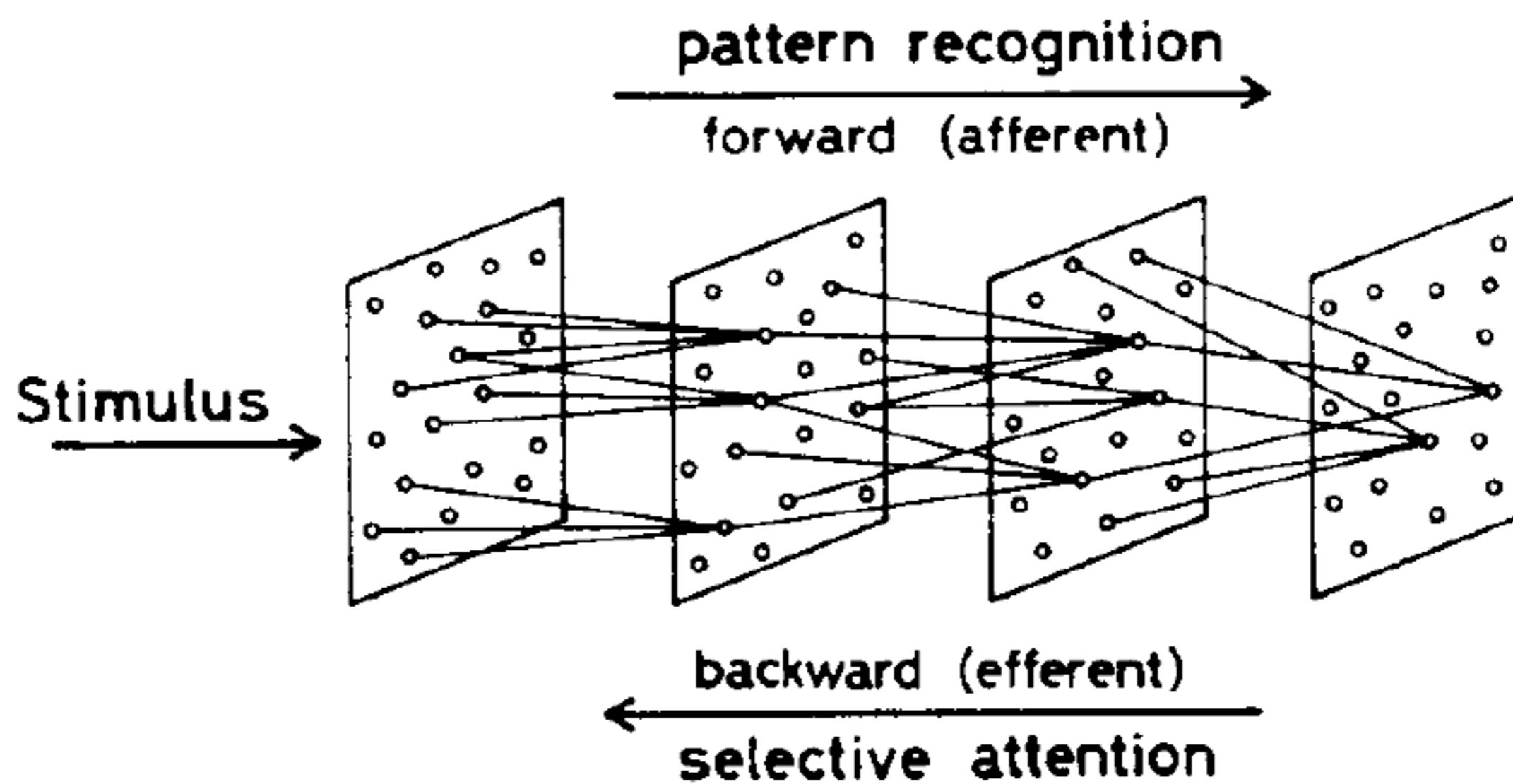


# Deep networks

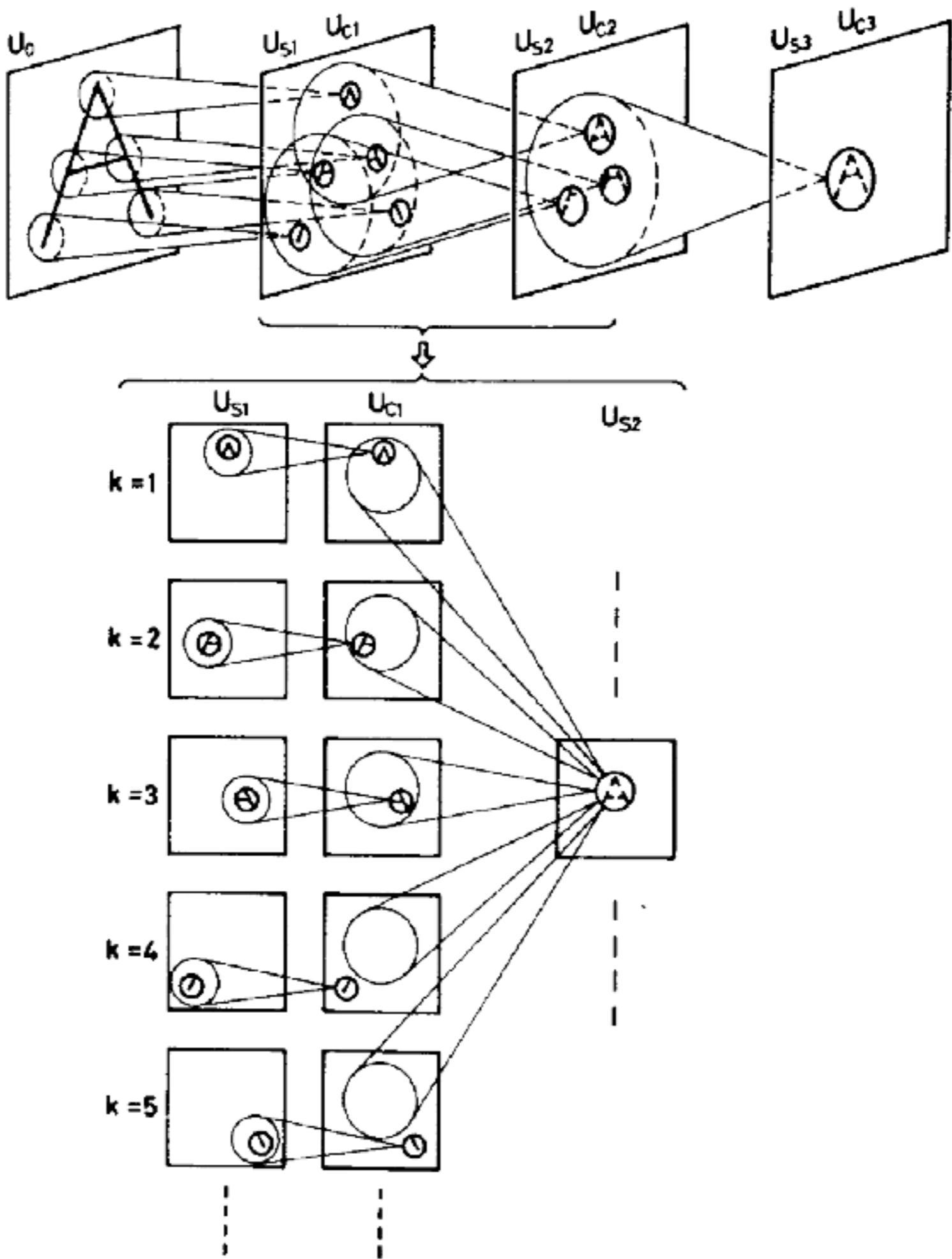
# Visual system







[Fukushima 1980]



[Fukushima 1980]

# LeNet + backpropagation (1988)

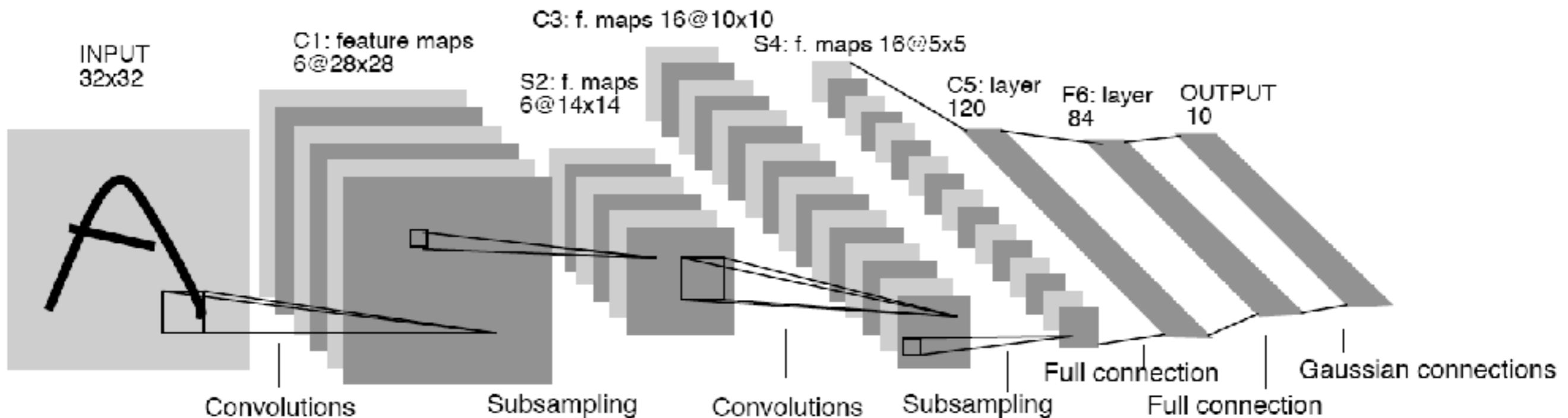


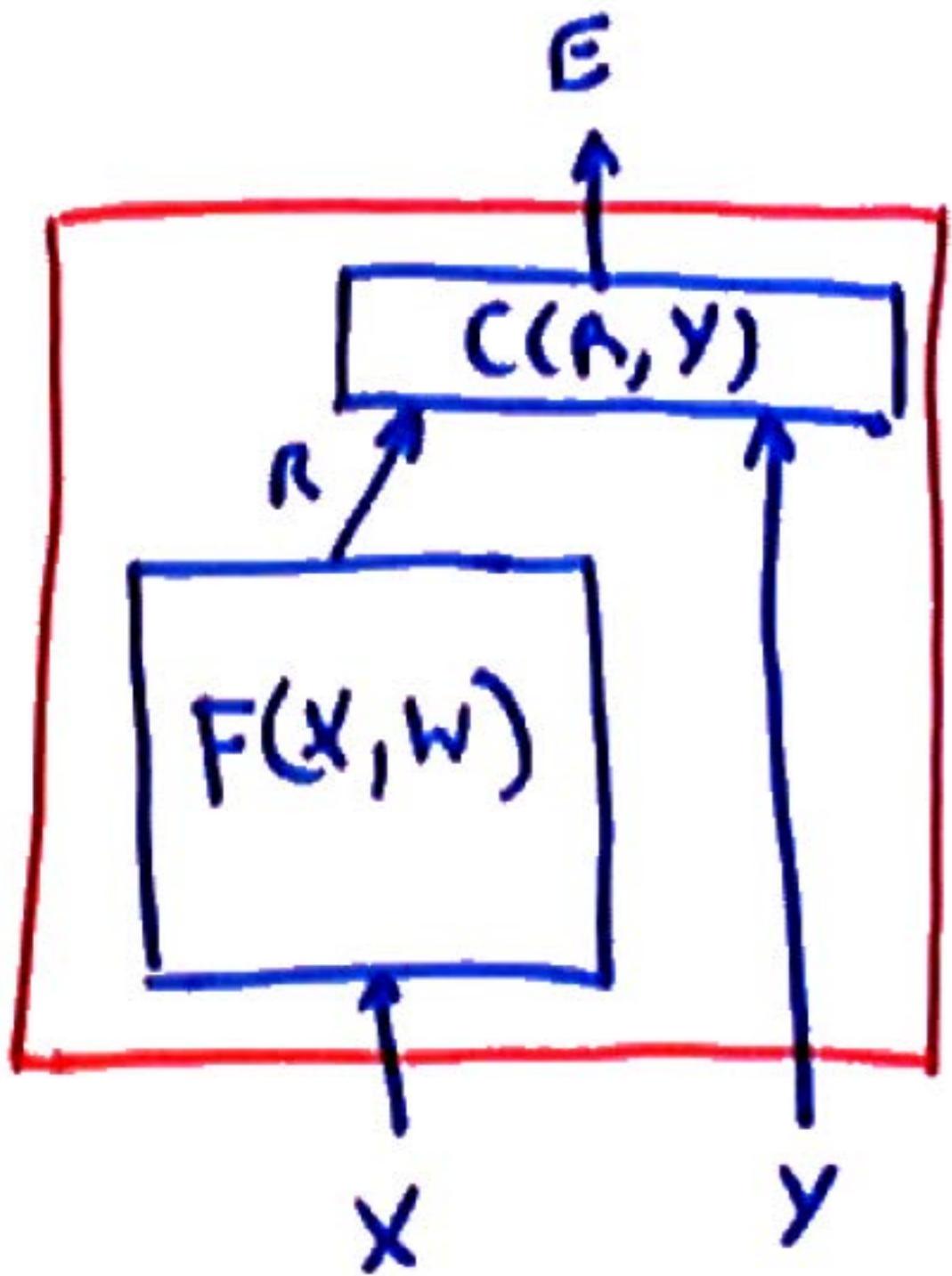
Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	X				X	X	X			X	X	X	X		X	X
1	X	X				X	X	X			X	X	X	X		X
2	X	X	X				X	X	X			X		X	X	X
3		X	X	X			X	X	X	X		X		X	X	X
4			X	X	X			X	X	X	X		X	X		X
5				X	X	X			X	X	X	X		X	X	X

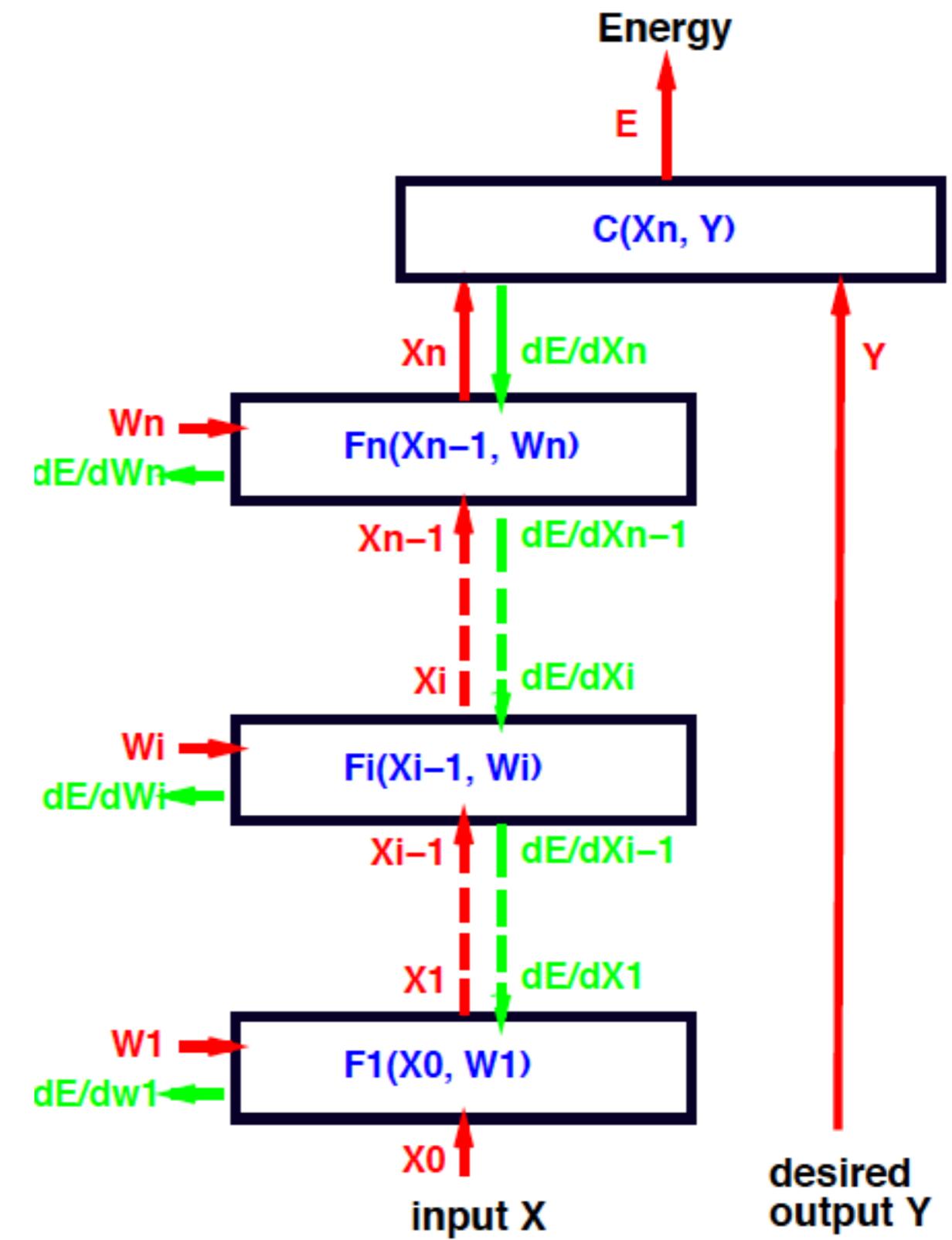
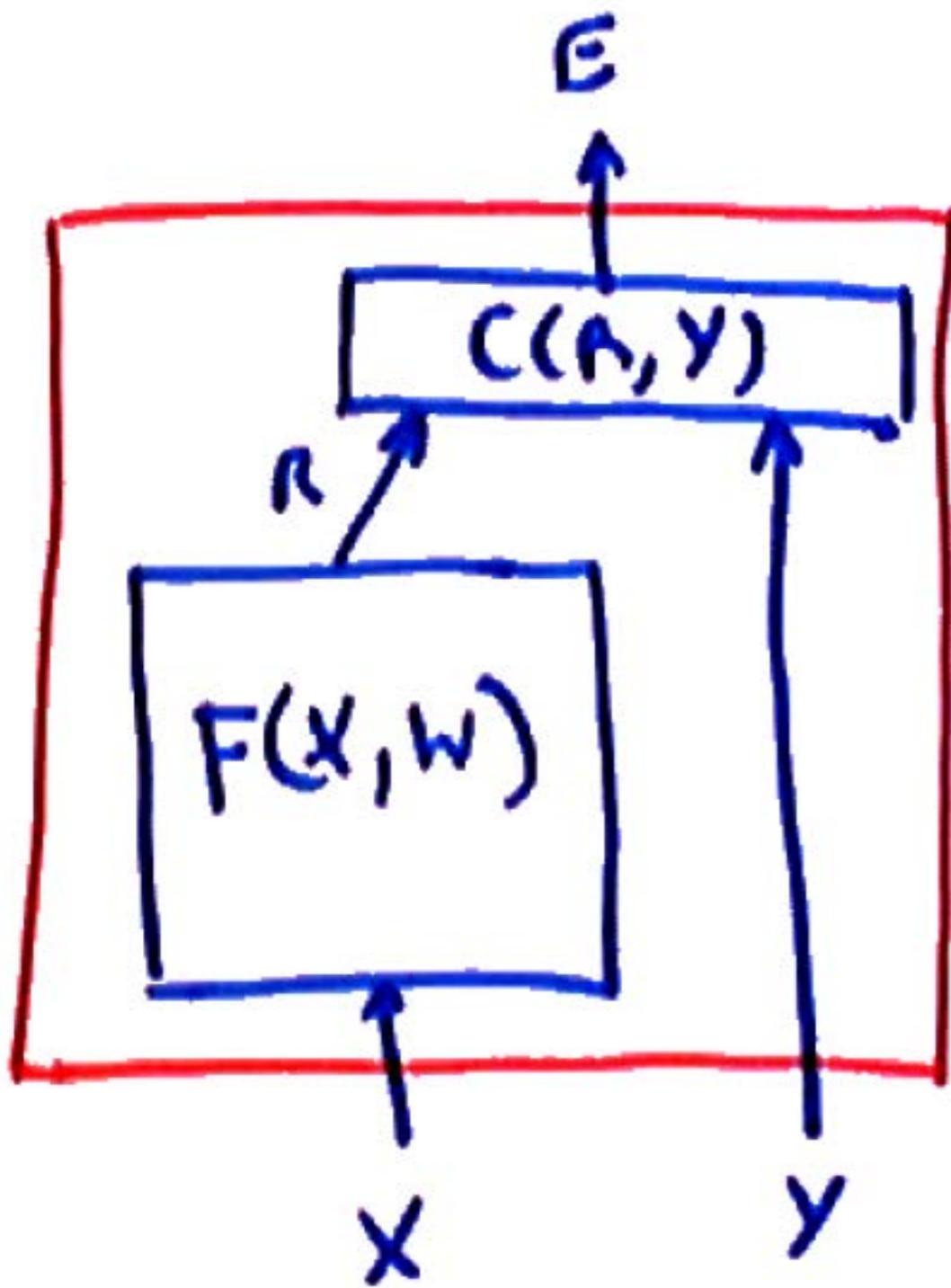
TABLE I  
EACH COLUMN INDICATES WHICH FEATURE MAP IN S2 ARE COMBINED  
BY THE UNITS IN A PARTICULAR FEATURE MAP OF C3.

[LeCun et al. 1988]

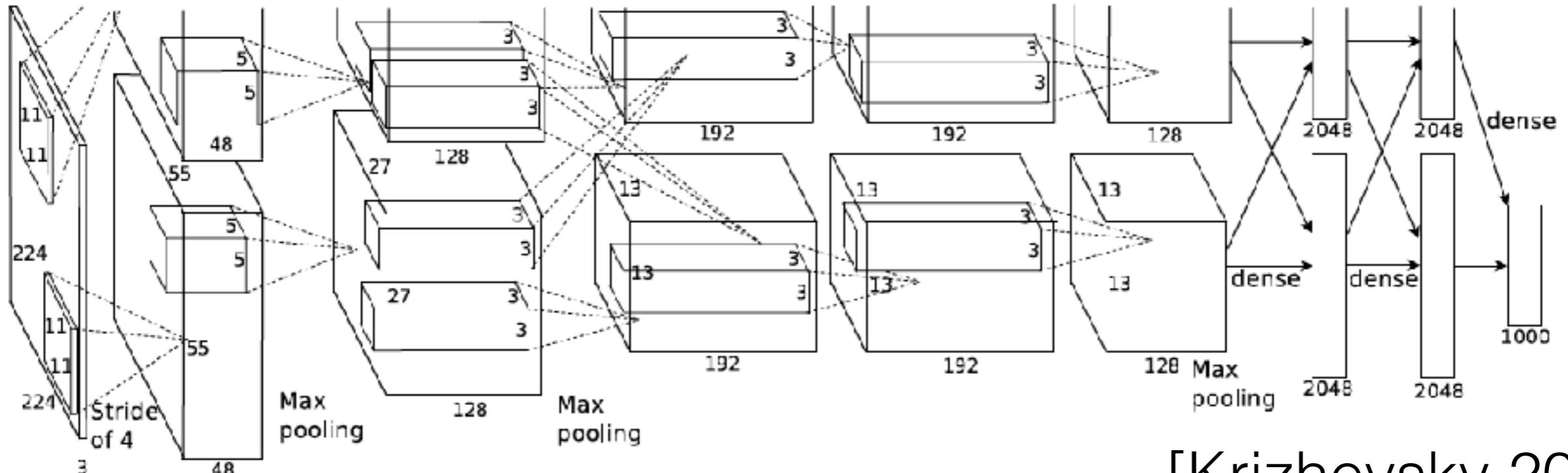
# Backpropagation



# Backpropagation



# AlexNet 2012



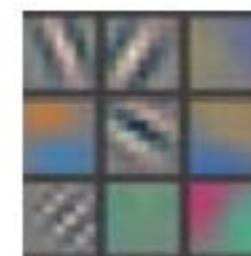
[Krizhevsky 2012]

**mite****container ship****motor scooter****leopard**

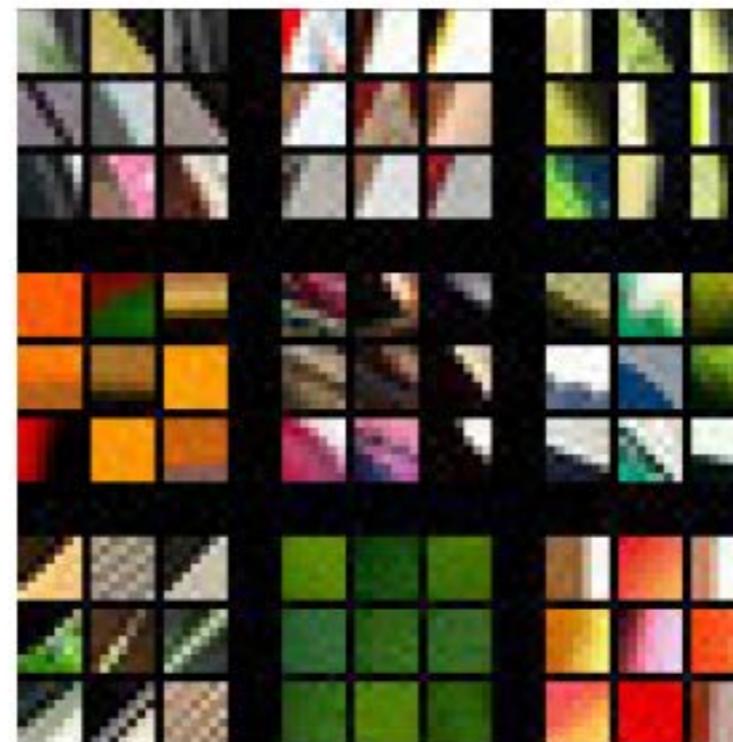
<b>mite</b>	<b>container ship</b>	<b>motor scooter</b>	<b>leopard</b>
black widow	lifeboat	go-kart	jaguar
cockroach	amphibian	moped	cheetah
tick	fireboat	bumper car	snow leopard
starfish	drilling platform	golfcart	Egyptian cat

**grille****mushroom****cherry****Madagascar cat**

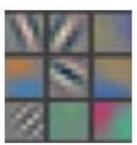
<b>convertible</b>	<b>agaric</b>	<b>dalmatian</b>	<b>squirrel monkey</b>
grille	mushroom	grape	spider monkey
pickup	jelly fungus	elderberry	titi
beach wagon	gill fungus	ffordshire bullterrier	indri
fire engine	dead-man's-fingers	currant	howler monkey



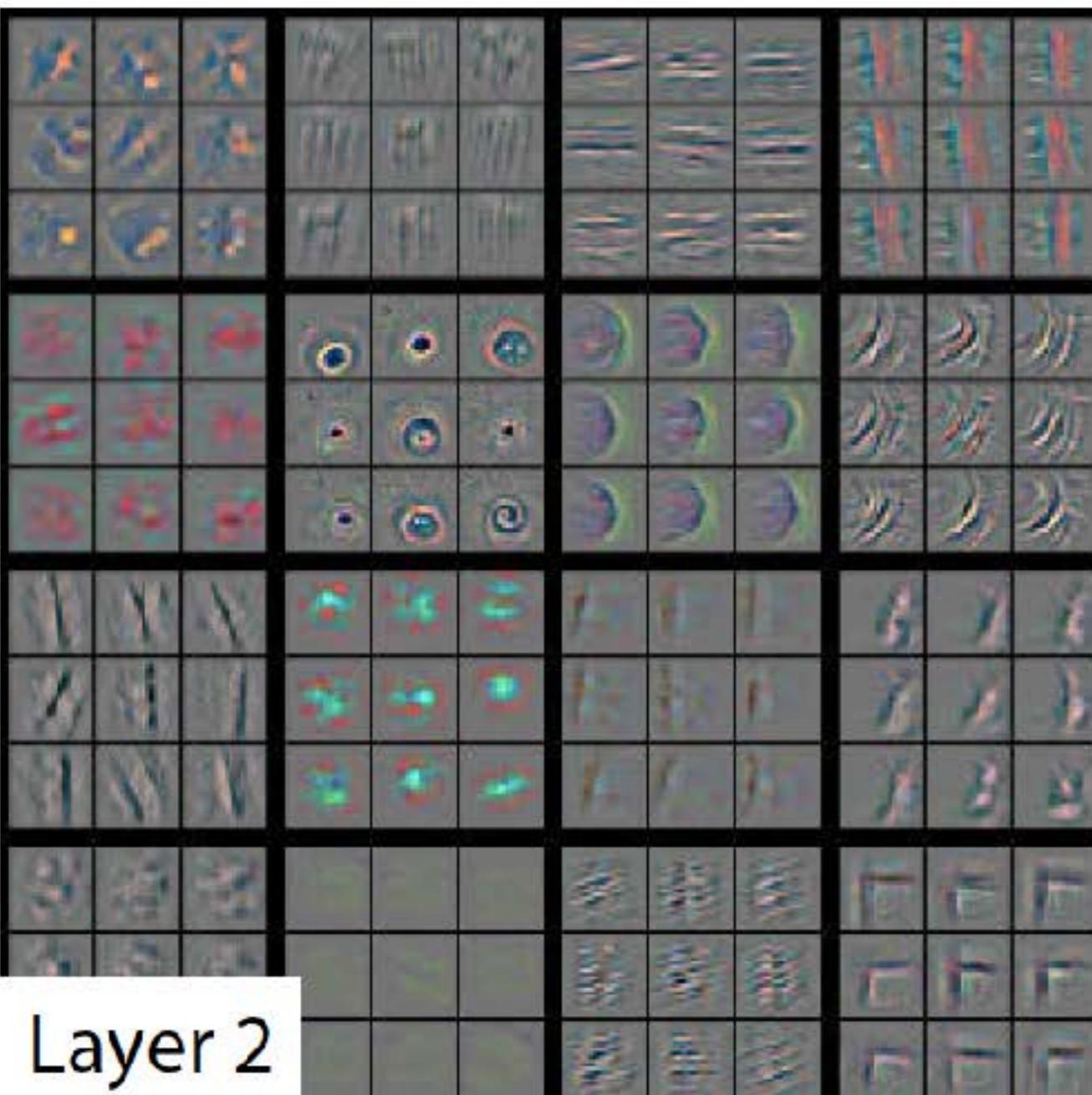
Layer 1



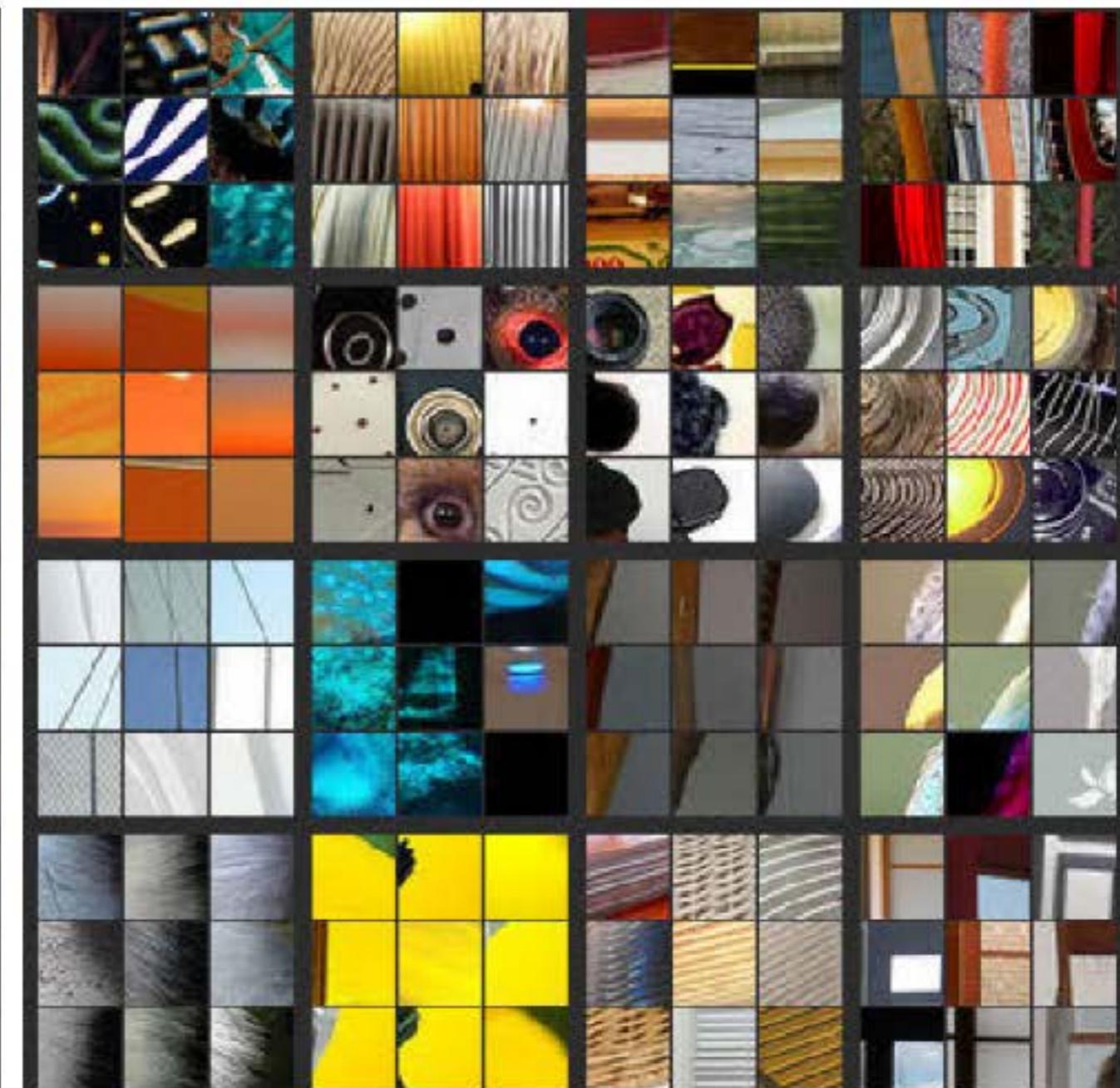
[Zeiler & Fergus 2014]

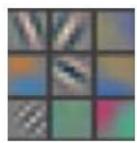


Layer 1

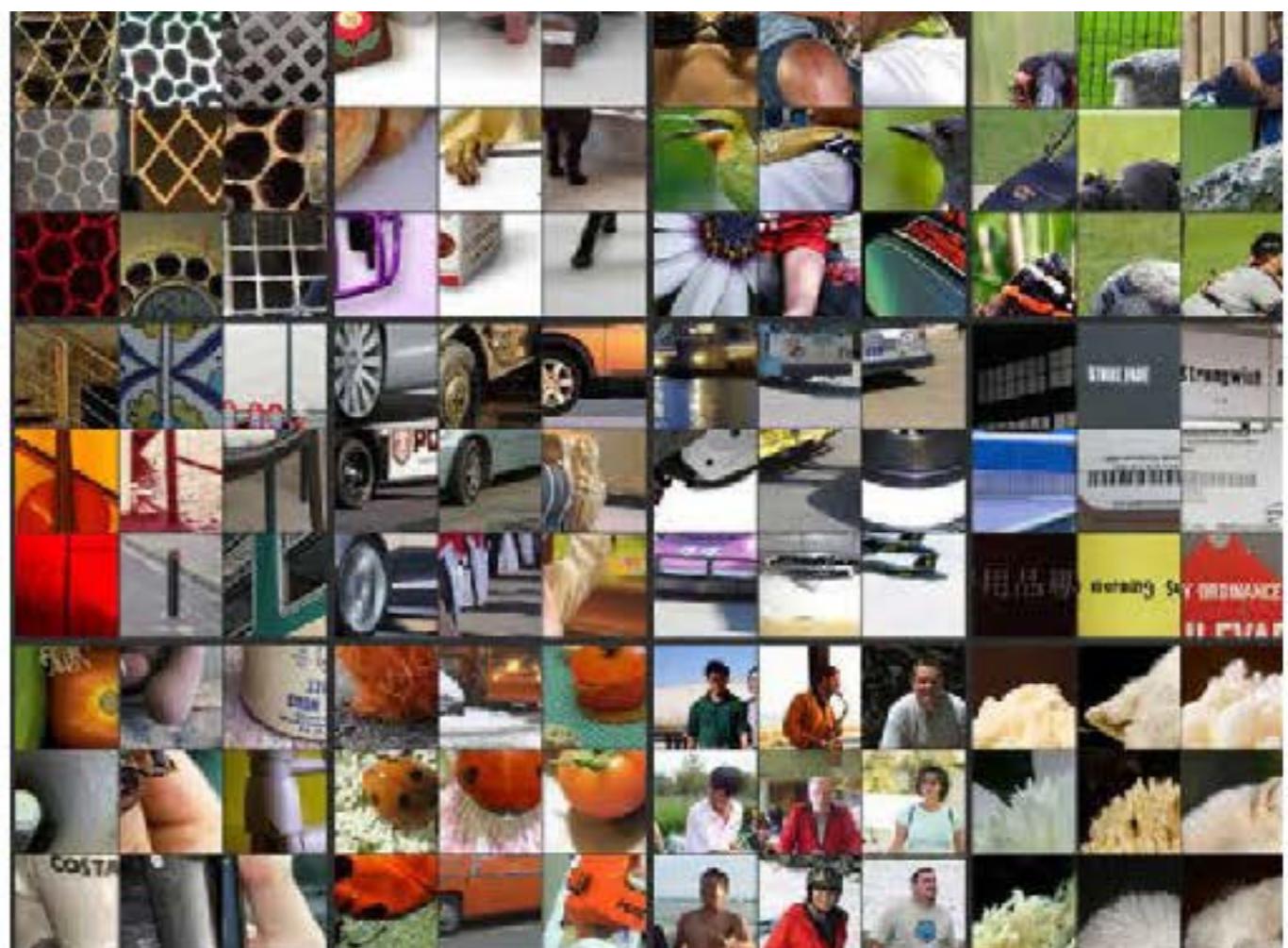
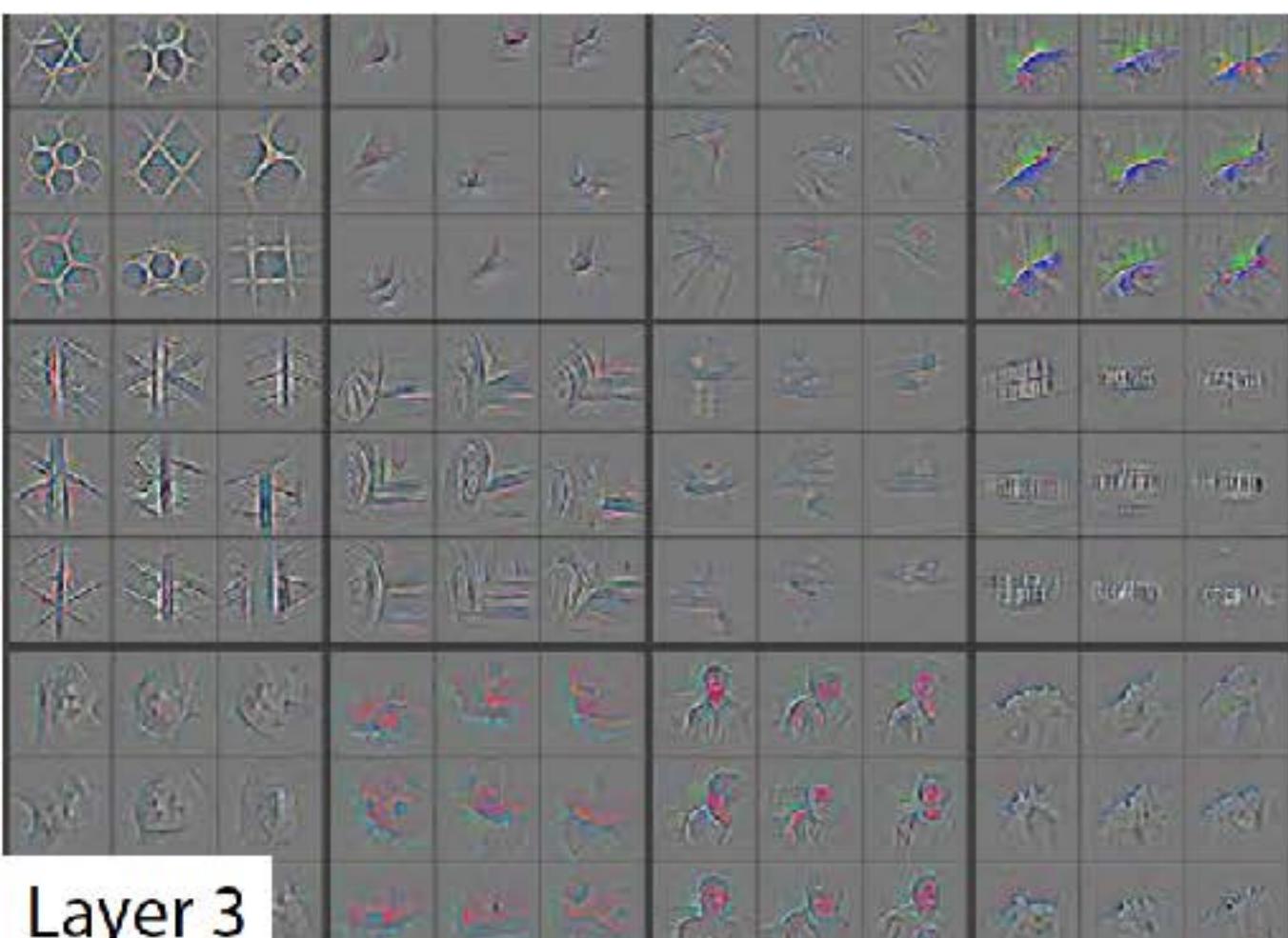
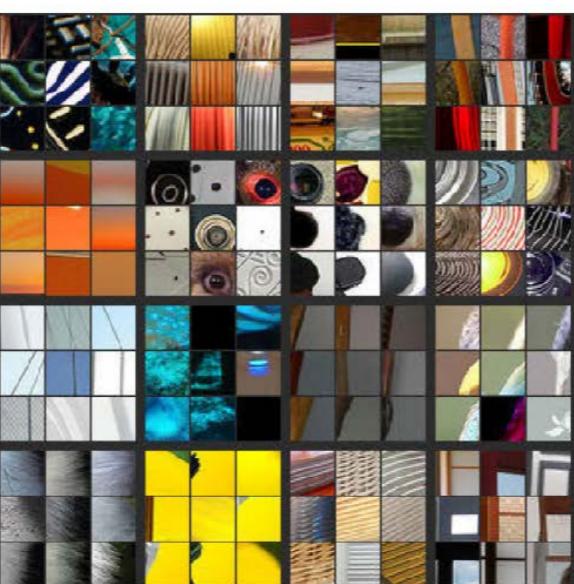
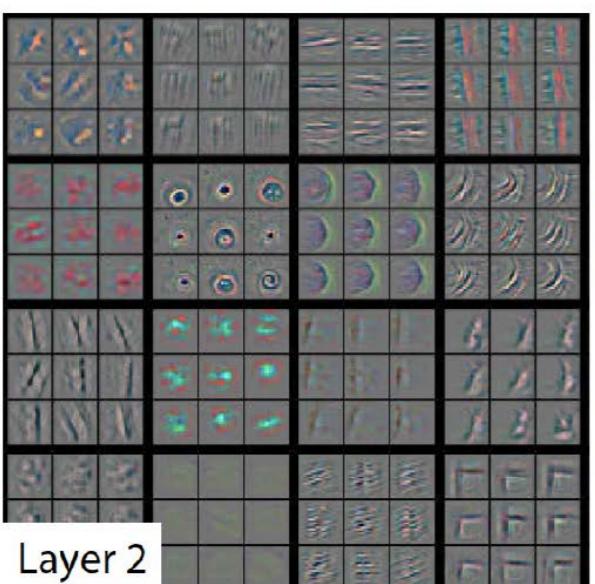


Layer 2

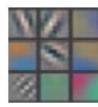




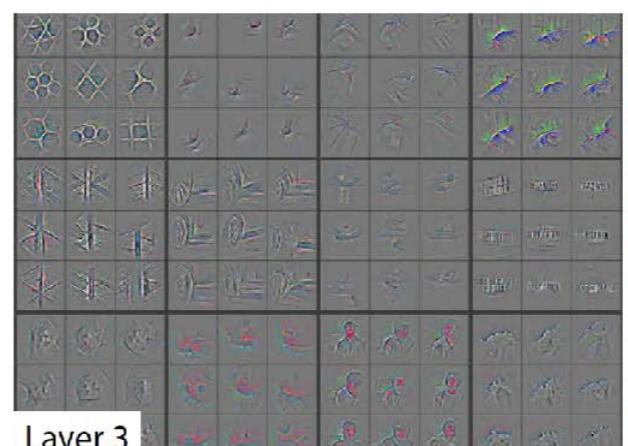
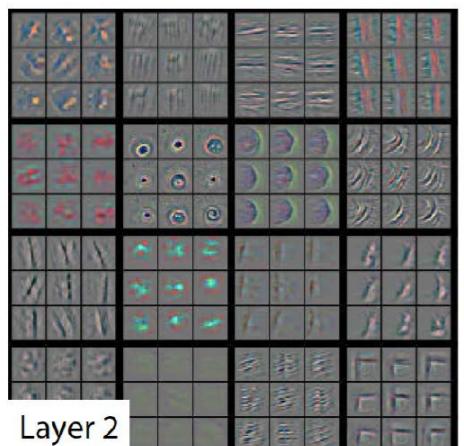
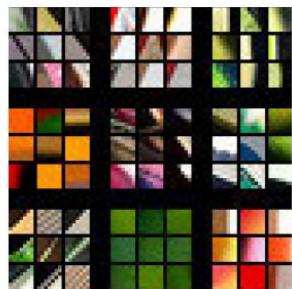
Layer 1



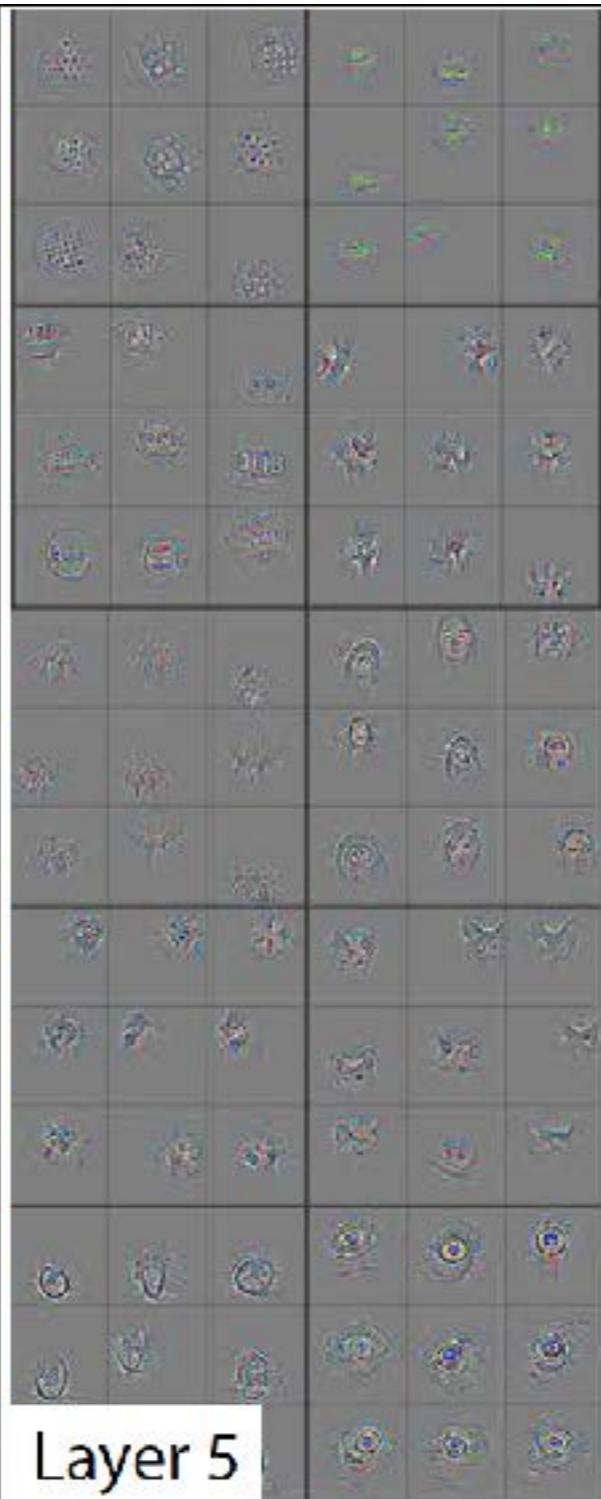
Layer 3



Layer 1



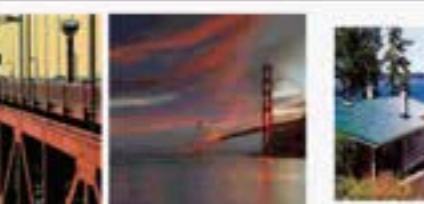
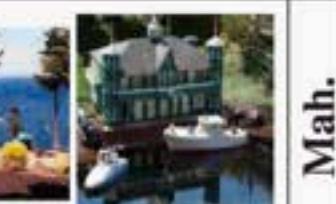
Layer 2



Layer 4

Layer 5

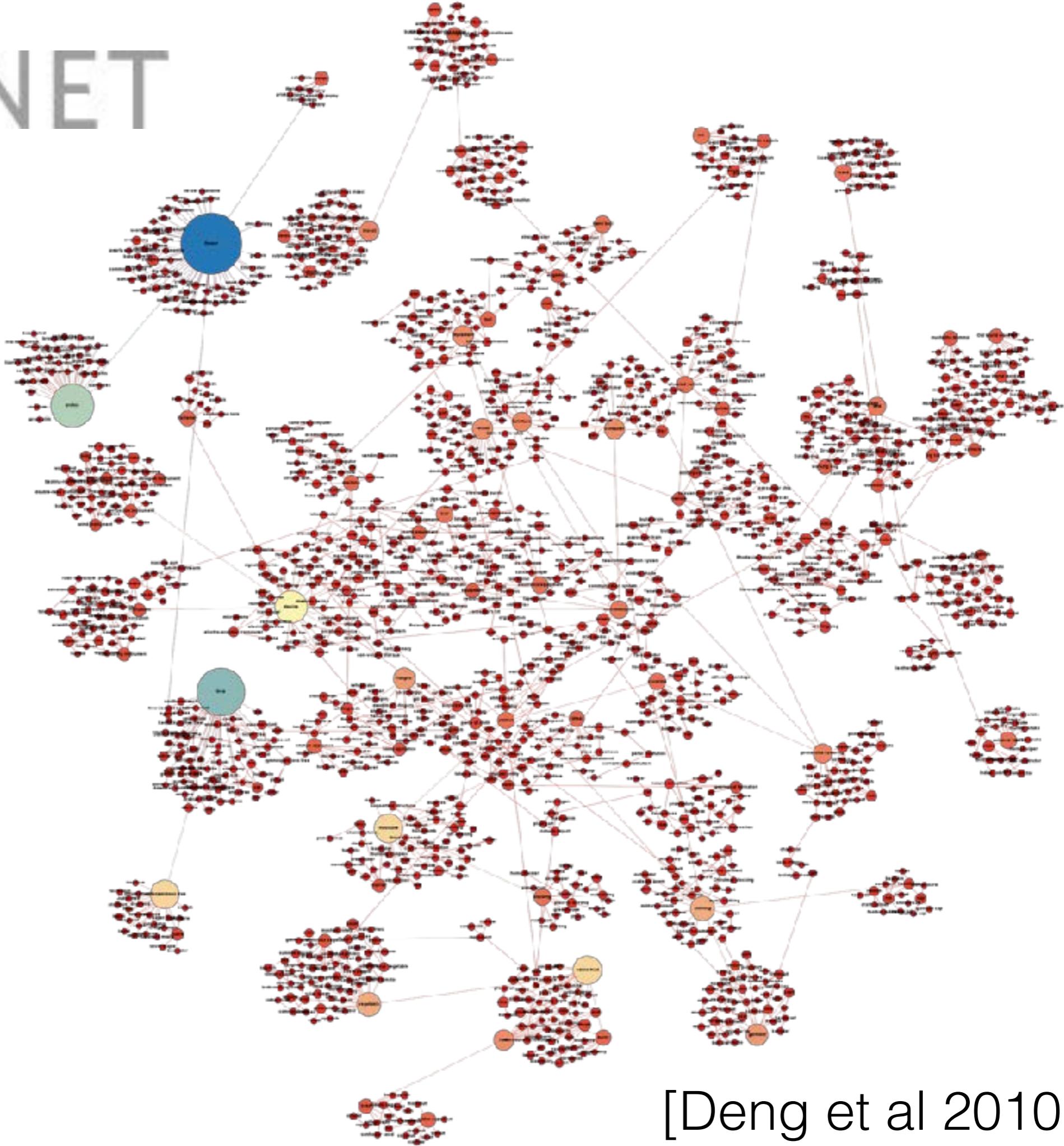
# Catalyst: ImageNet 1000

L2	shopping cart 1.07%	unicycle 0.84%	covered wagon 0.83%	garbage truck 0.79%	forklift 0.78%
Gondola L2 4.4% - Mah. 99.7%	 	 	 	 	 
	 dock 0.11%	 canoe 0.03%	 fishing rod 0.01%	 bridge 0.01%	 boathouse 0.01%
	 crane 0.87%	 stupa 0.83%	 roller coaster 0.79%	 bell cote 0.78%	 flagpole 0.75%
Palm L2 6.4% - Mah. 98.1%	 	 cabbage tree 0.81%	 pine 0.30%	 pandanus 0.14%	 iron tree 0.07%
					 logwood 0.06%

~1K images for 1K categories = ~1M images

[Deng et al 2010]

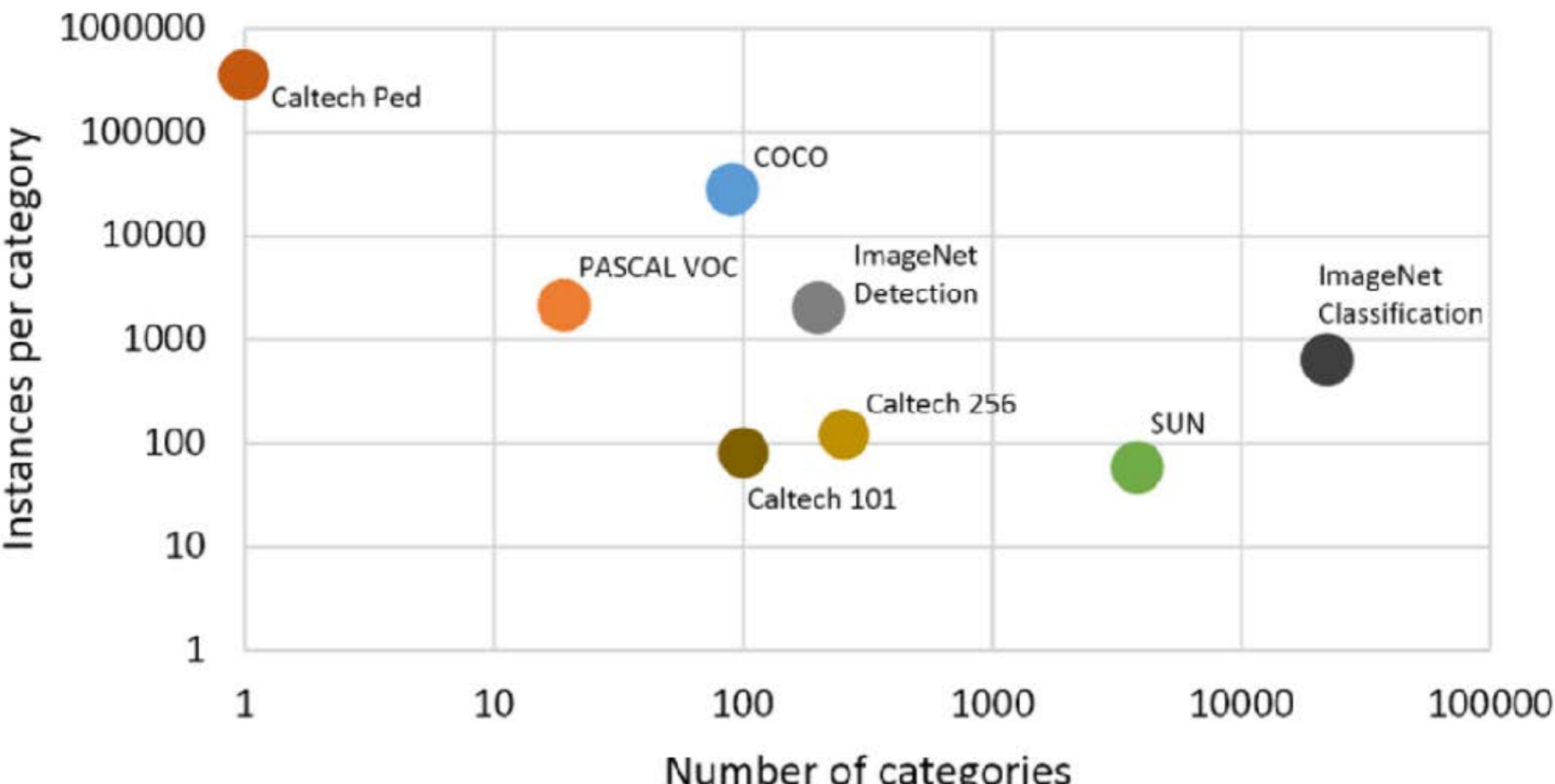
# IMAGENET



[Deng et al 2010]

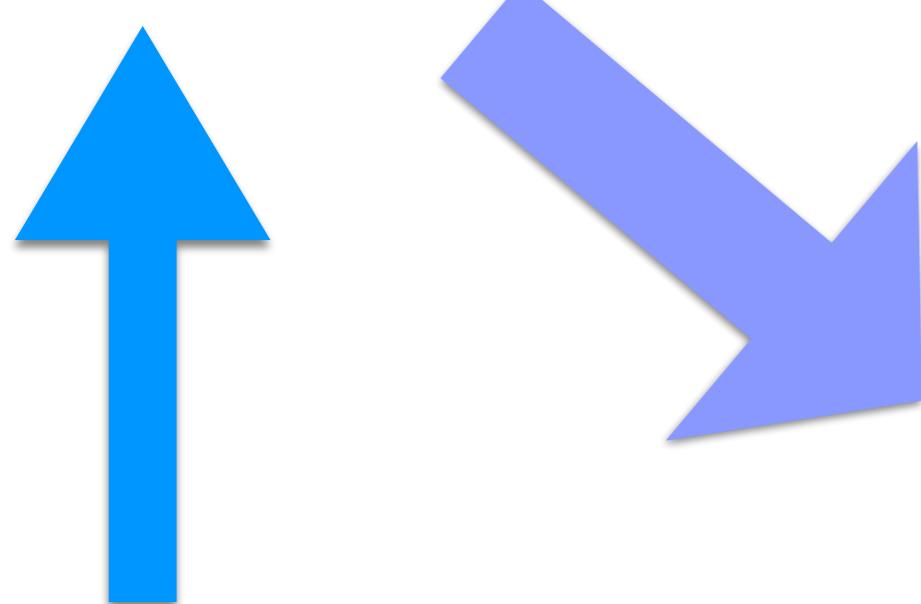
# Large annotated datasets

Number of categories vs. number of instances



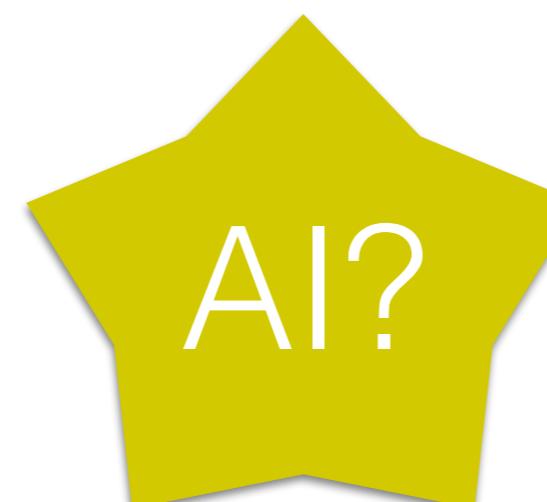
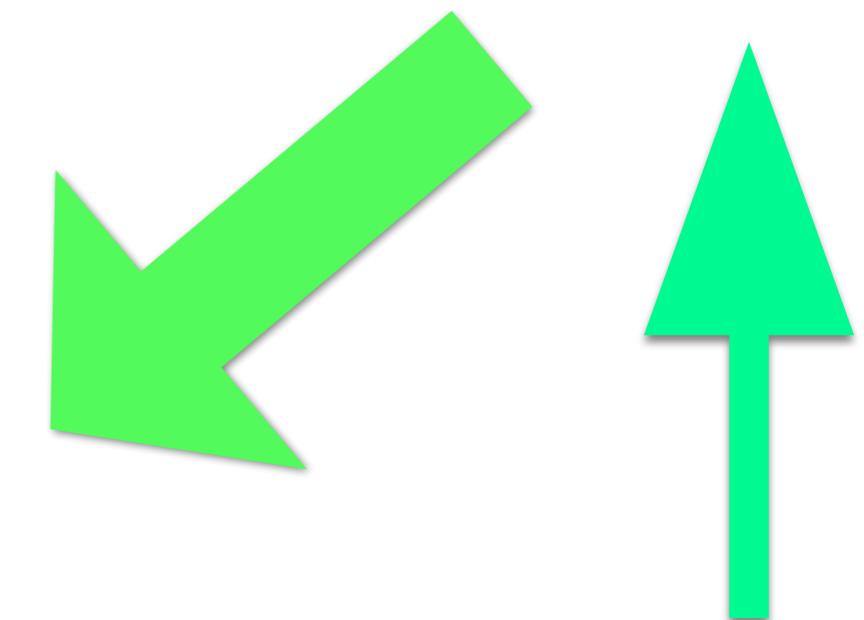
[Lin et al. 2015]

Deep networks  
1980-1990

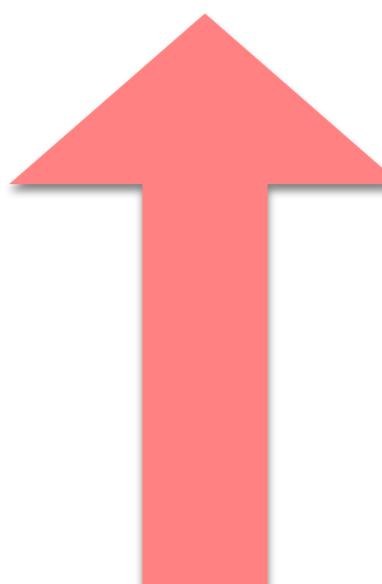


Neuroscience  
1960-1990

Large annotated datasets  
2004-2010

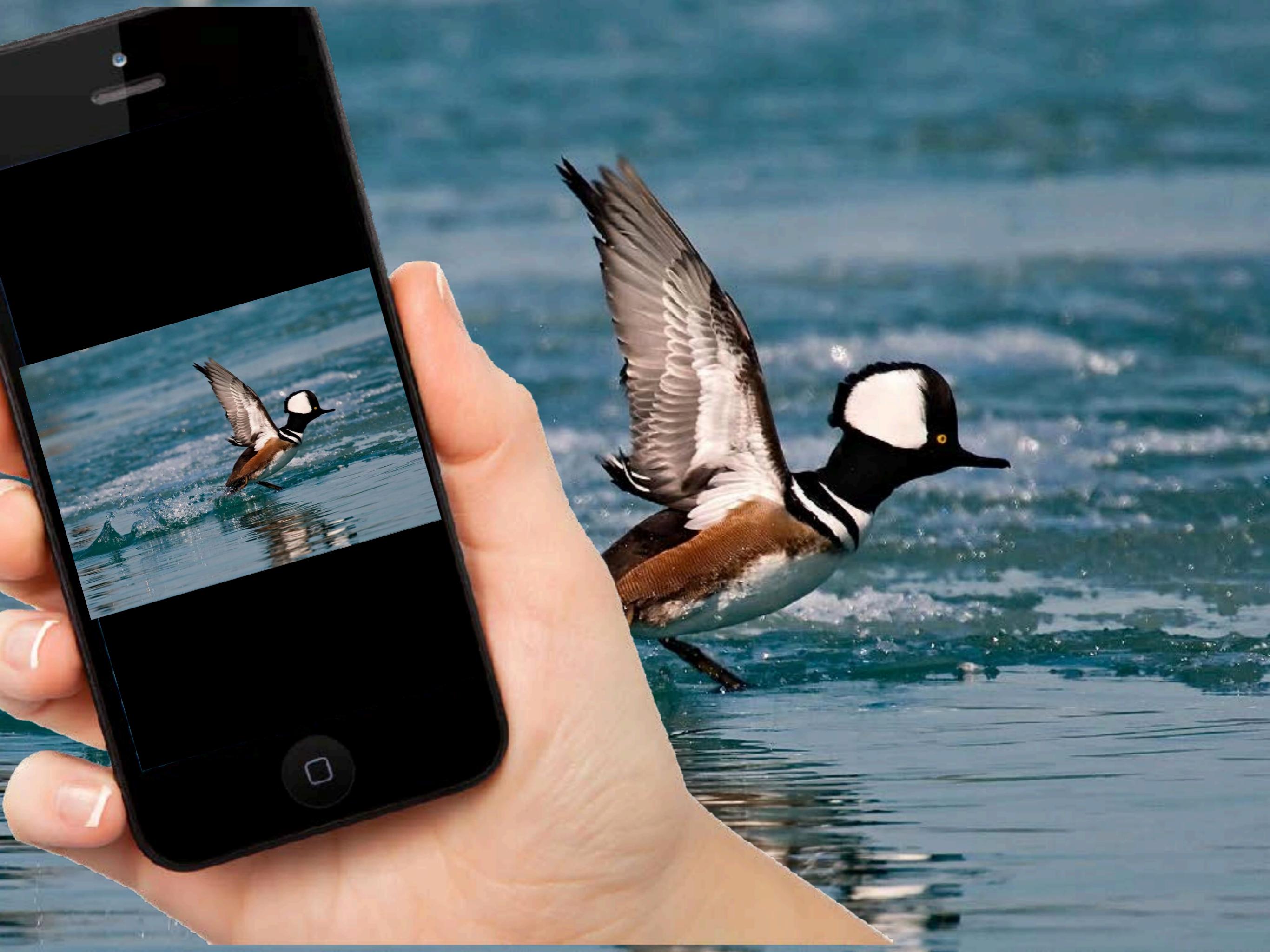


Google,  
Flickr, AMT  
2003-2008



Moore's law, GPUs  
1960-2015

success stories



550 species  
of N. American  
birds

Verizon 16:59 90%

< Edit Detail List Home

Hooded Merganser

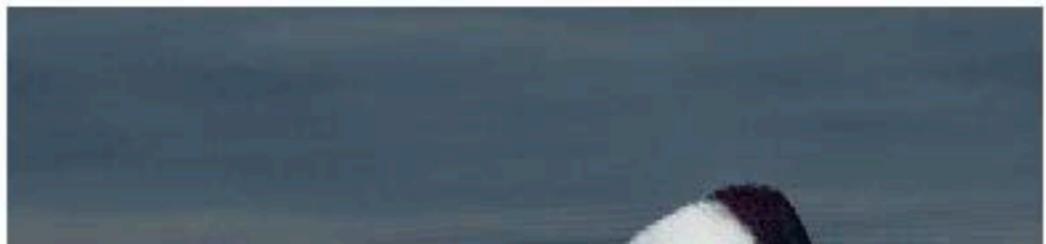


BREEDING MALE

Small duck; feeds by diving to catch mainly fish with thin, serrated bill. Breeding males have showy black and white crest, a coupl...

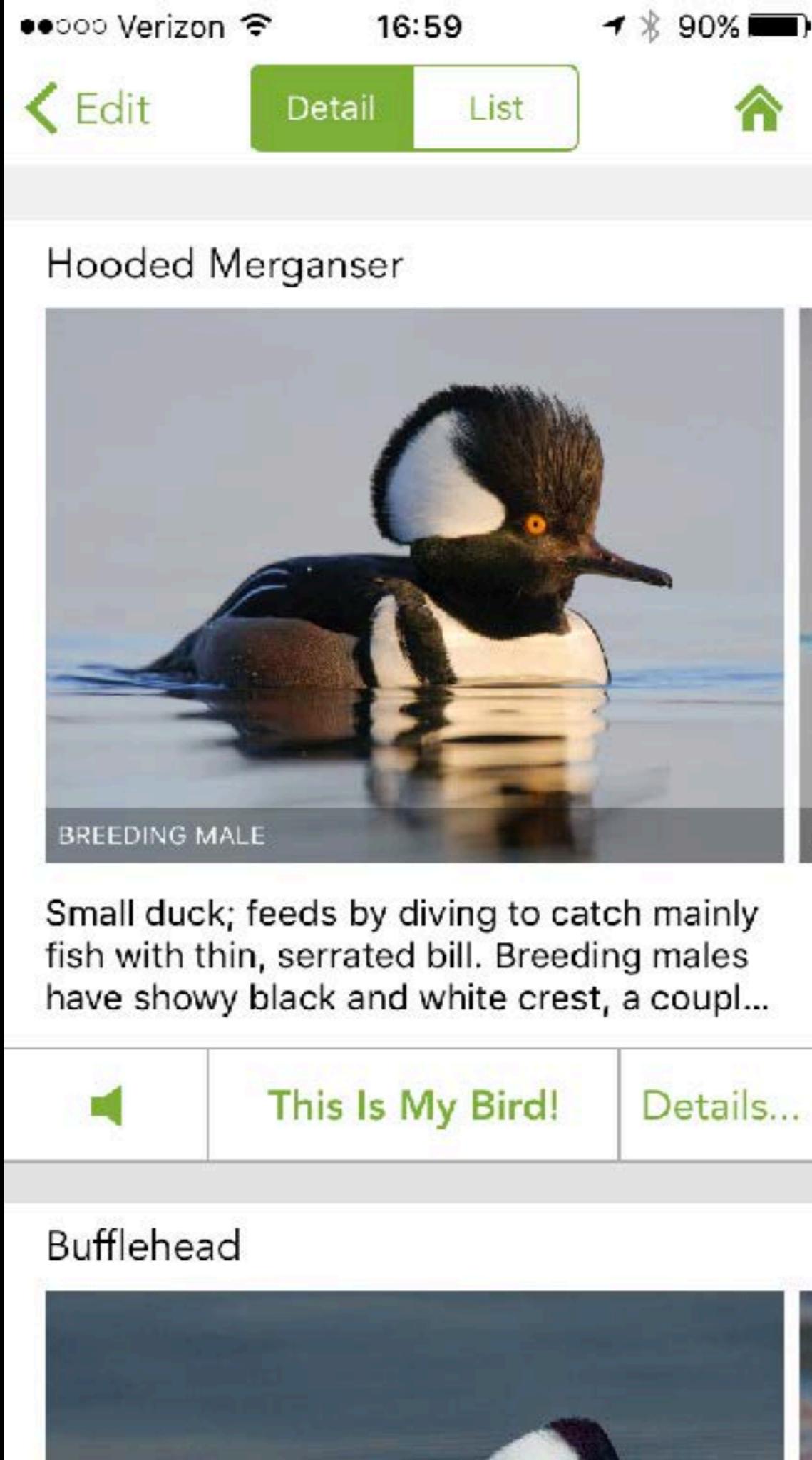
Speaker This Is My Bird! Details...

Bufflehead



[van Horn et al  
2014,2016]

550 species  
of N. American  
birds



App store:  
``Merlin Bird ID''

[van Horn et al  
2014,2016]

• SPARROWS are small brown-bodied birds with streaked backs and short conical beaks. Their food, mostly seeds except during the nesting season, is obtained on or near the ground. When not nesting, most are seen because they are sho-

these may  
be identi-  
fied by  
pevies are  
presented

# Sparrows

by one head each; the striking Lark Bunting, the Snow Bunting, the towhees, and the Olive Sparrow are omitted. Immatures of some species are much duller, especially those species with black or rufous on the head. Songs and chips of sparrows are often more easily distinguished than are their plumages. See pp. 328-345 for further details.

## STREAKED BREASTS



Vesper  
p. 332  
(white outer tail)

Song  
p. 342

Lincoln's  
p. 342

Savannah  
p. 328



Le Conte's  
p. 330



Sharp-tailed  
p. 330



Henslow's  
p. 328  
(rufous wing)



Baird's  
p. 328



Purple Finch  
for comparison  
p. 316



Seaside  
p. 330



Fox  
p. 342  
(rufous rump)

Sage  
p. 332

## UNSTREAKED BREASTS



Dark-eyed Junco  
p. 334  
(white outer tail)

Black-chinned  
p. 335

Black-throated  
p. 332  
(white outer tail)

Lapland Longspur  
p. 344

White-crowned  
p. 310

White-throated  
p. 340

Golden-crowned  
p. 340

Harris'  
p. 340

American Tree  
p. 338

Field  
p. 338

Chipping  
p. 338  
(gray rump)

Swamp  
p. 342

Brewer's  
p. 338

Clay-colored  
p. 338  
(brown rump)

Grasshopper  
p. 328

Rufous-crowned  
p. 336

Lark  
p. 332  
(white tail fringe)

Rufous-winged  
p. 336

Cassin's  
p. 336

Bartram's  
p. 336

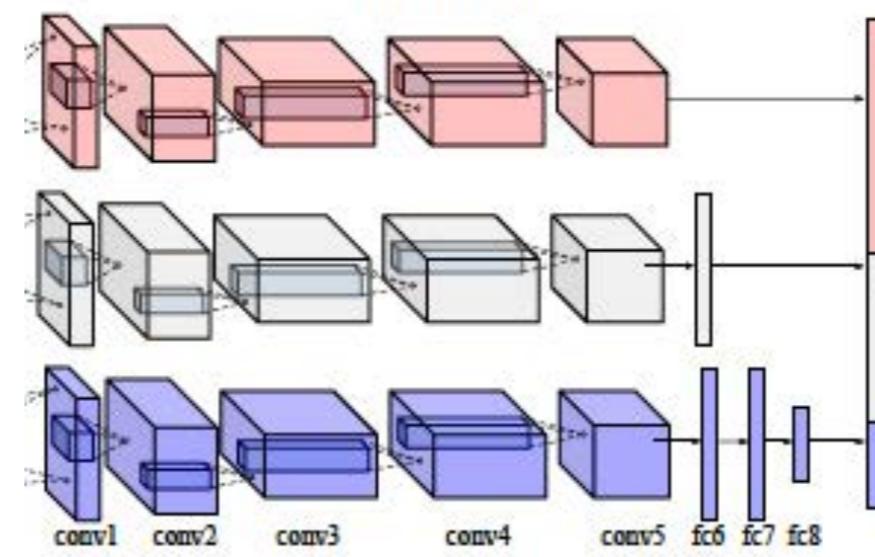
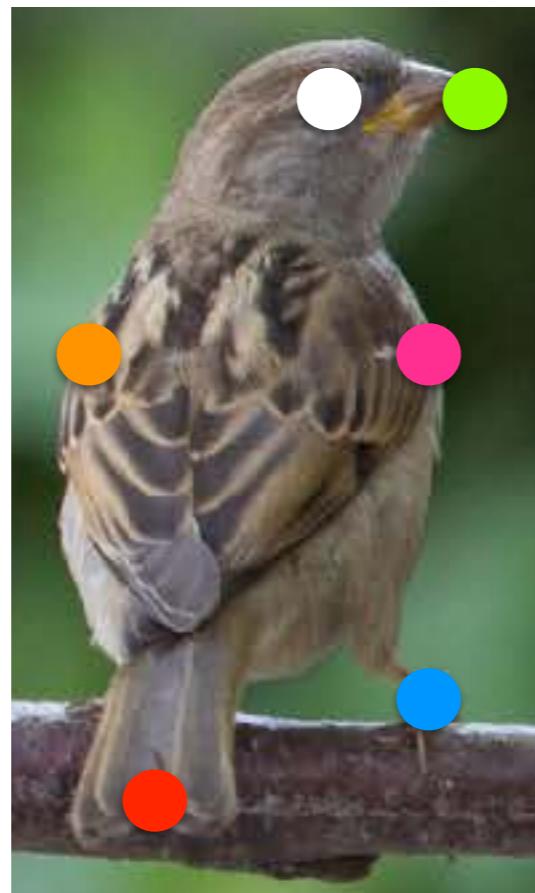
# Invariance to pose and background

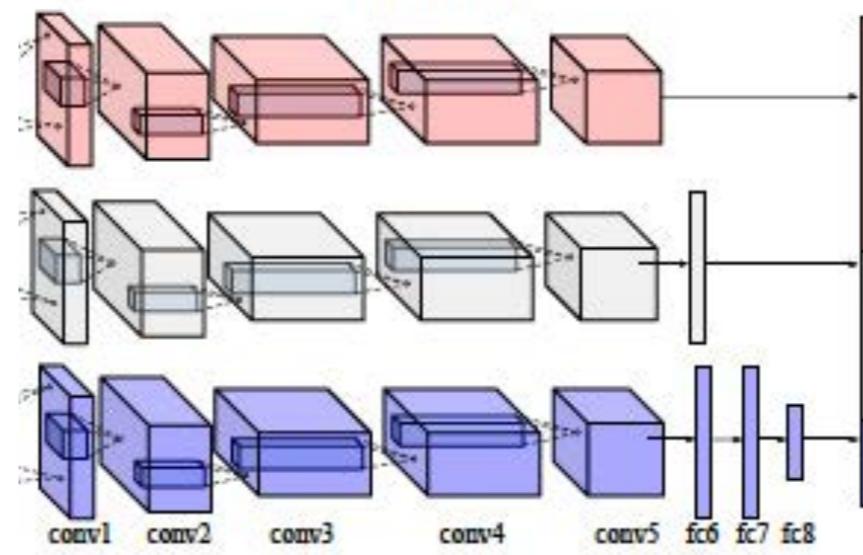










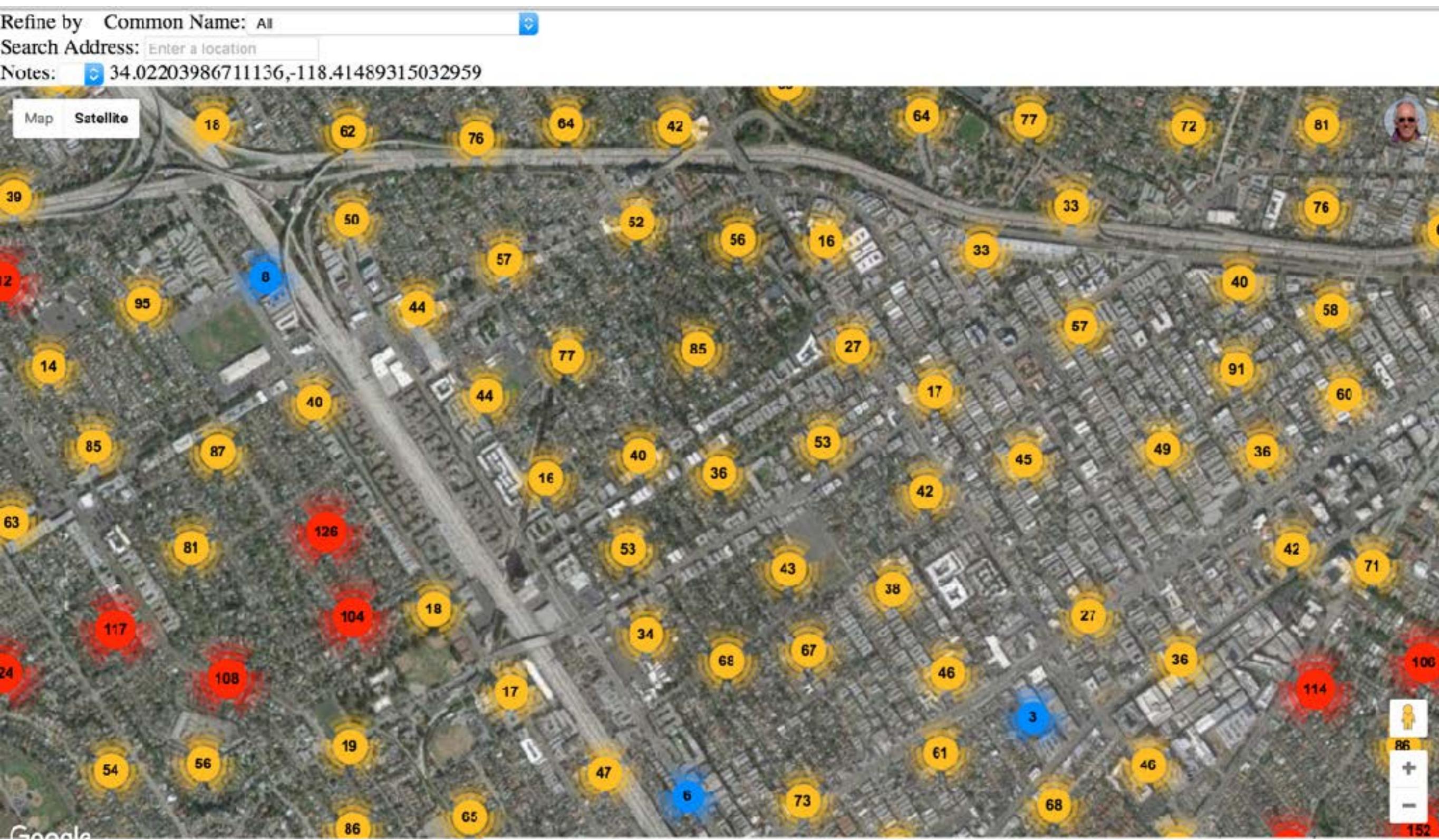


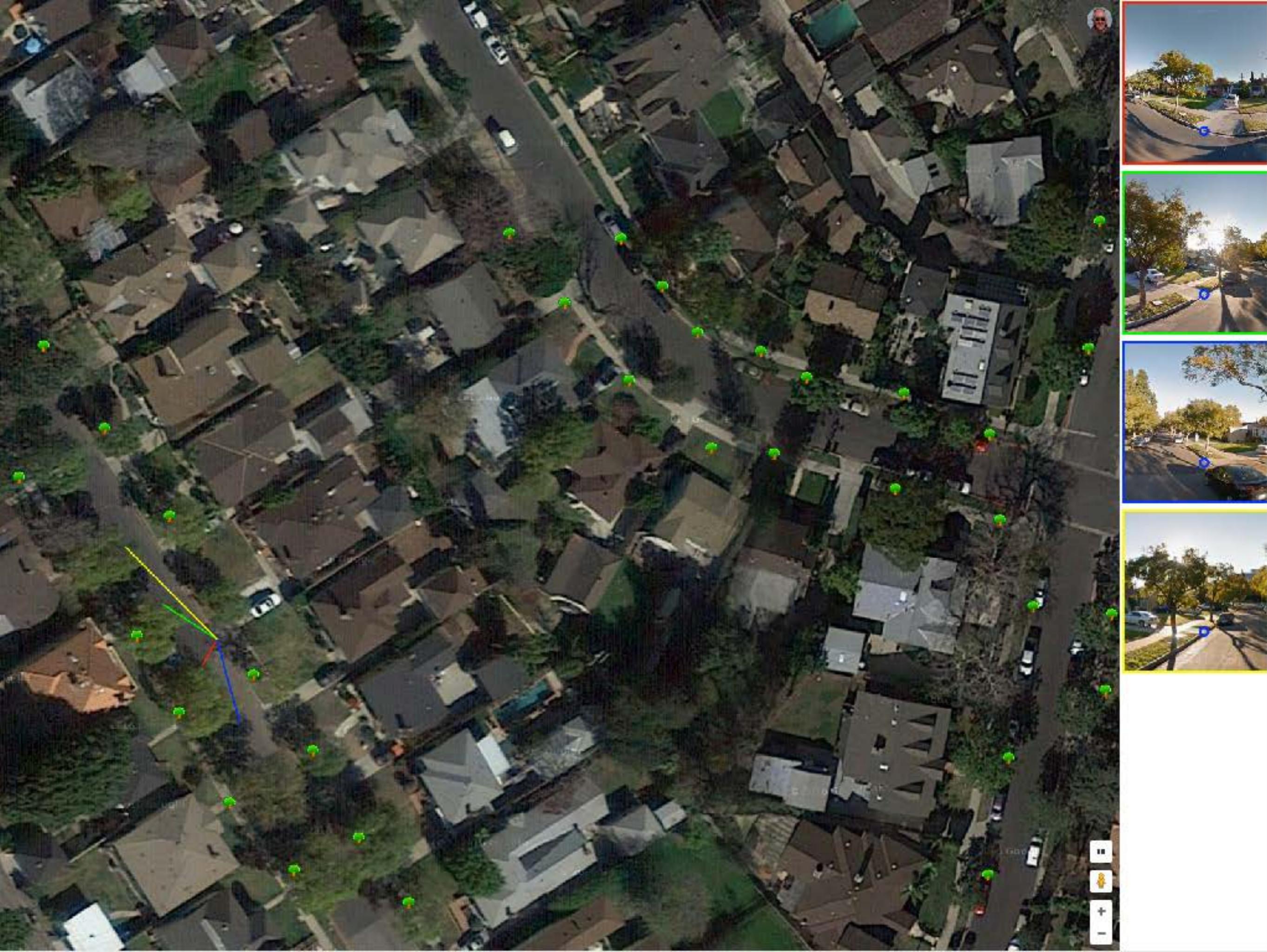
‘house  
sparrow’



[van Horn et al. 2016]

# Los Angeles = 1M trees





ID: 548

Location: 34.032781281453966, 118.42148597780546

Species List: Jacaranda(2.01), Chinese sweetgum(-0.384), Dracaena palmi(1.01), Magnolia(-1.018), Rough Shell Macadamia(-1.095)

Trunk Diameter: 15

Detection Score: 0.222

Street View: N8cocimR:T1J8C4o4V3cxUO

Localization Error (m): N/A



[Update Note](#)



Registree Caltech



b B

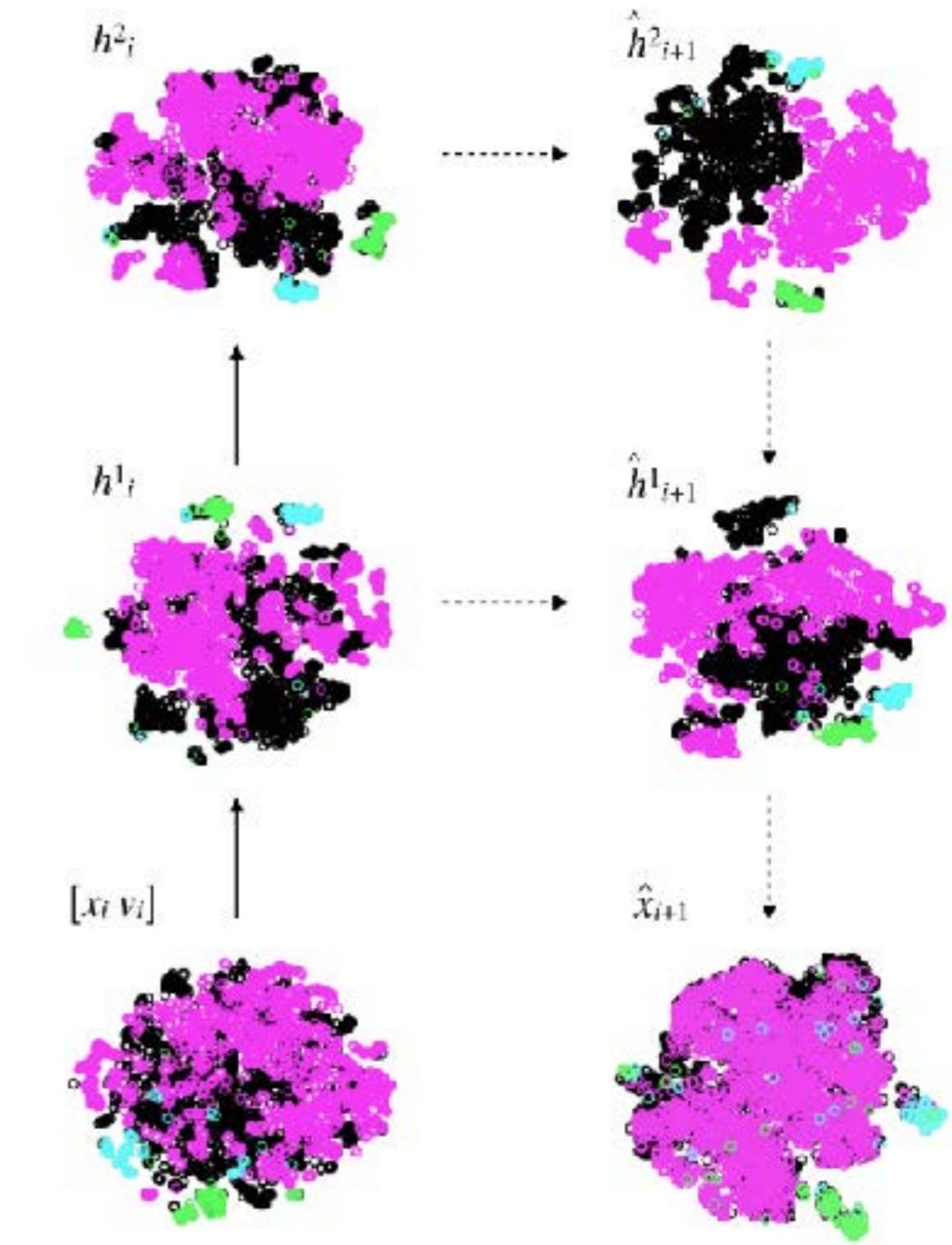
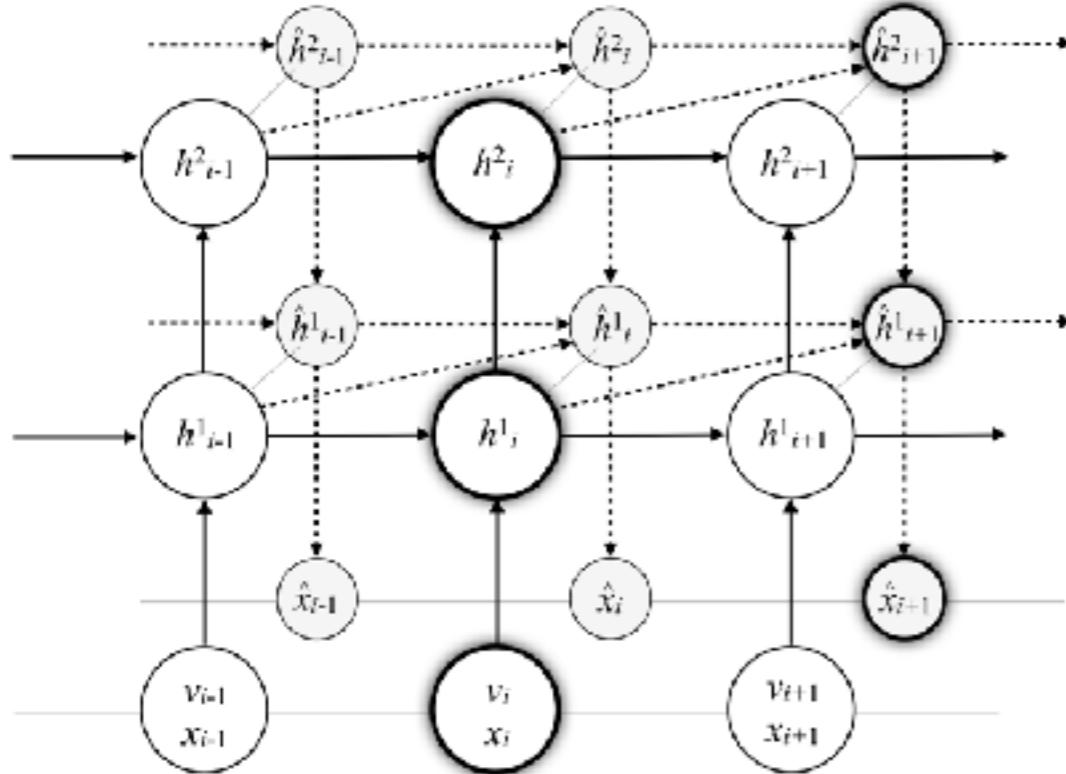


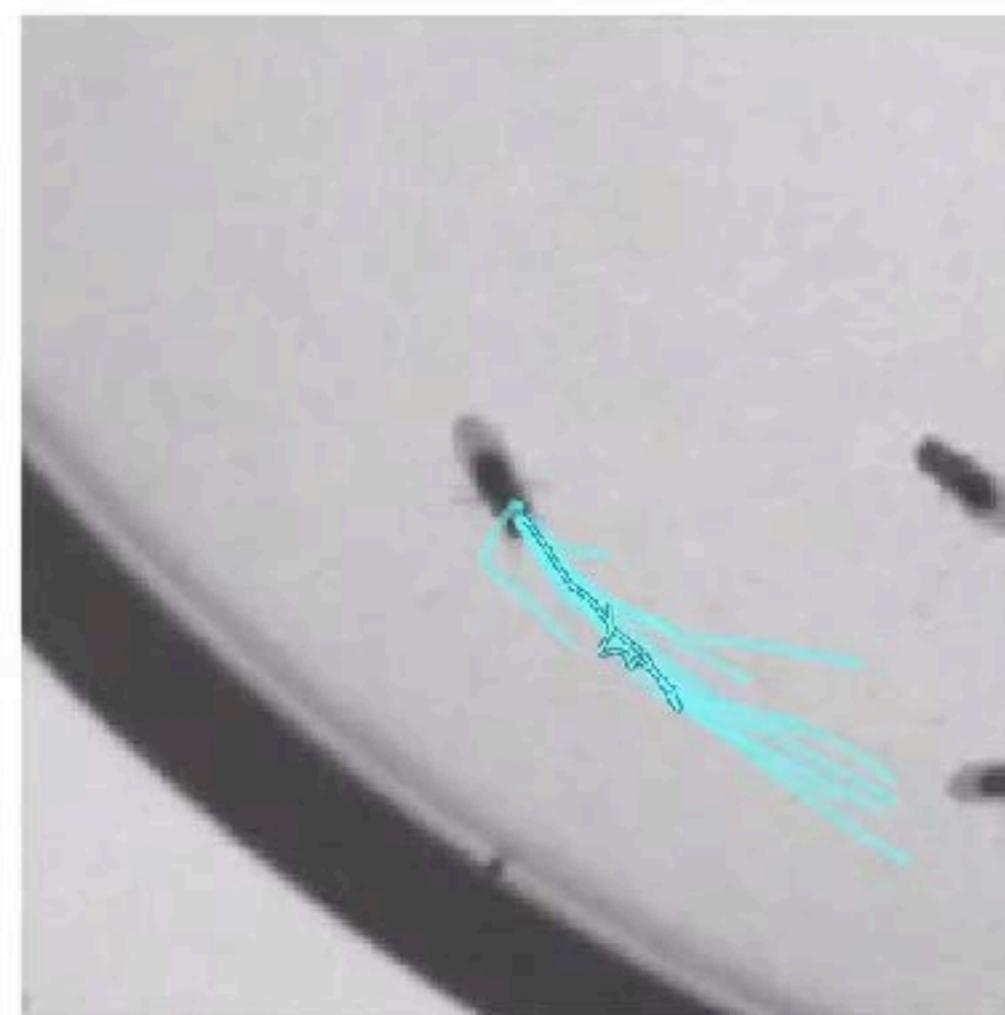
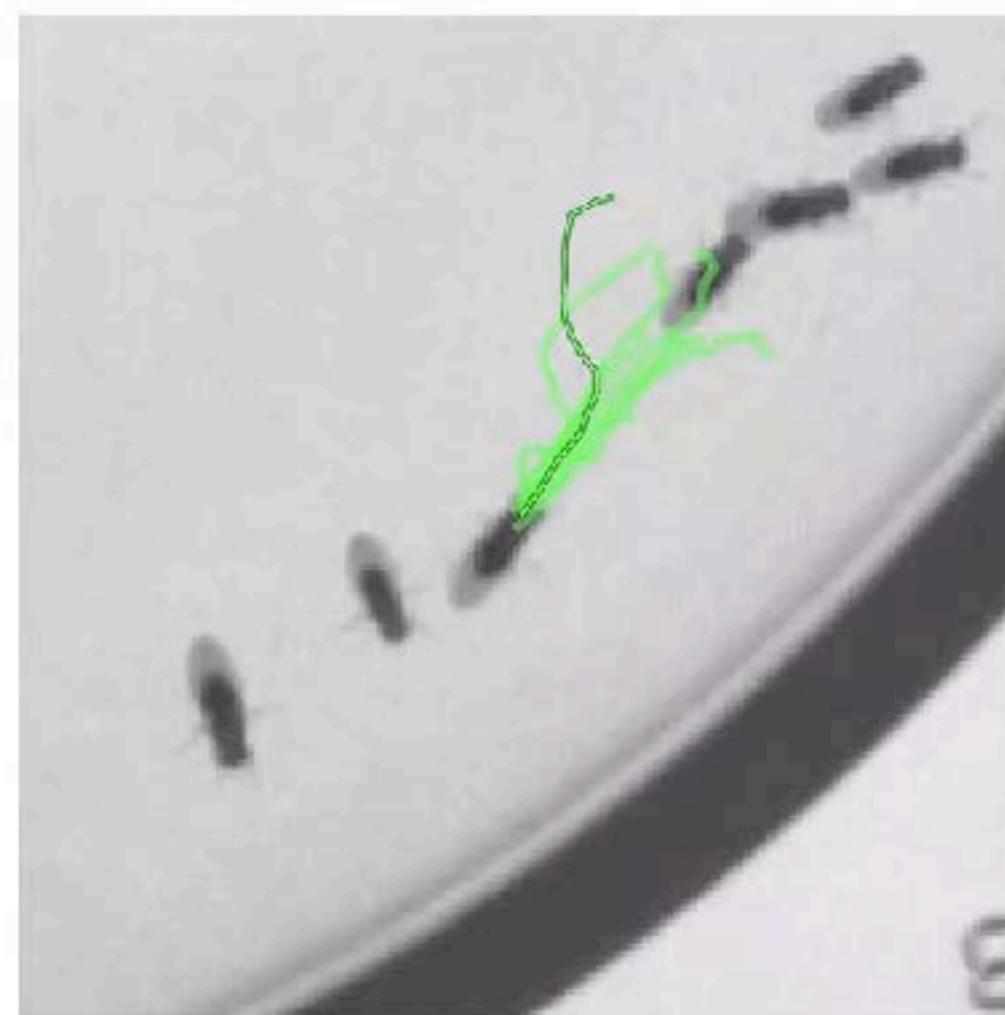
# Discovery: FlyBowl

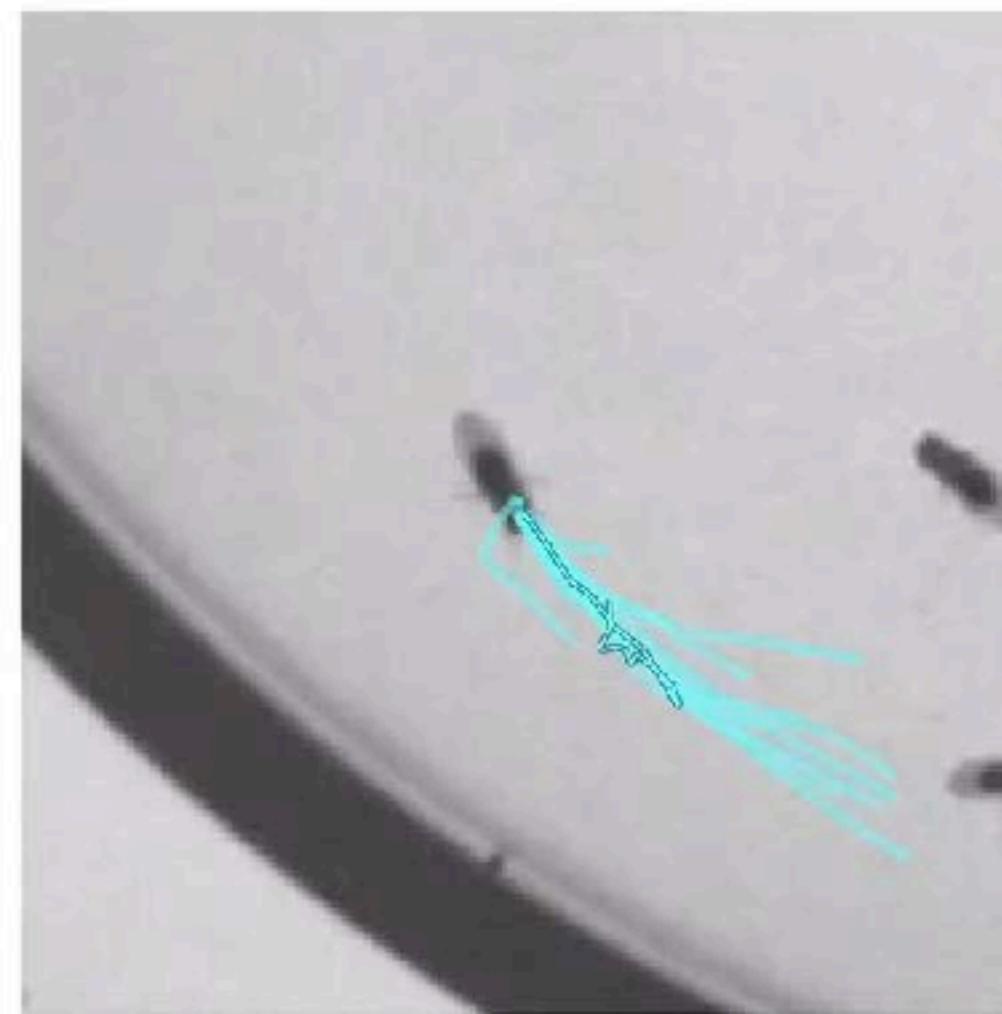
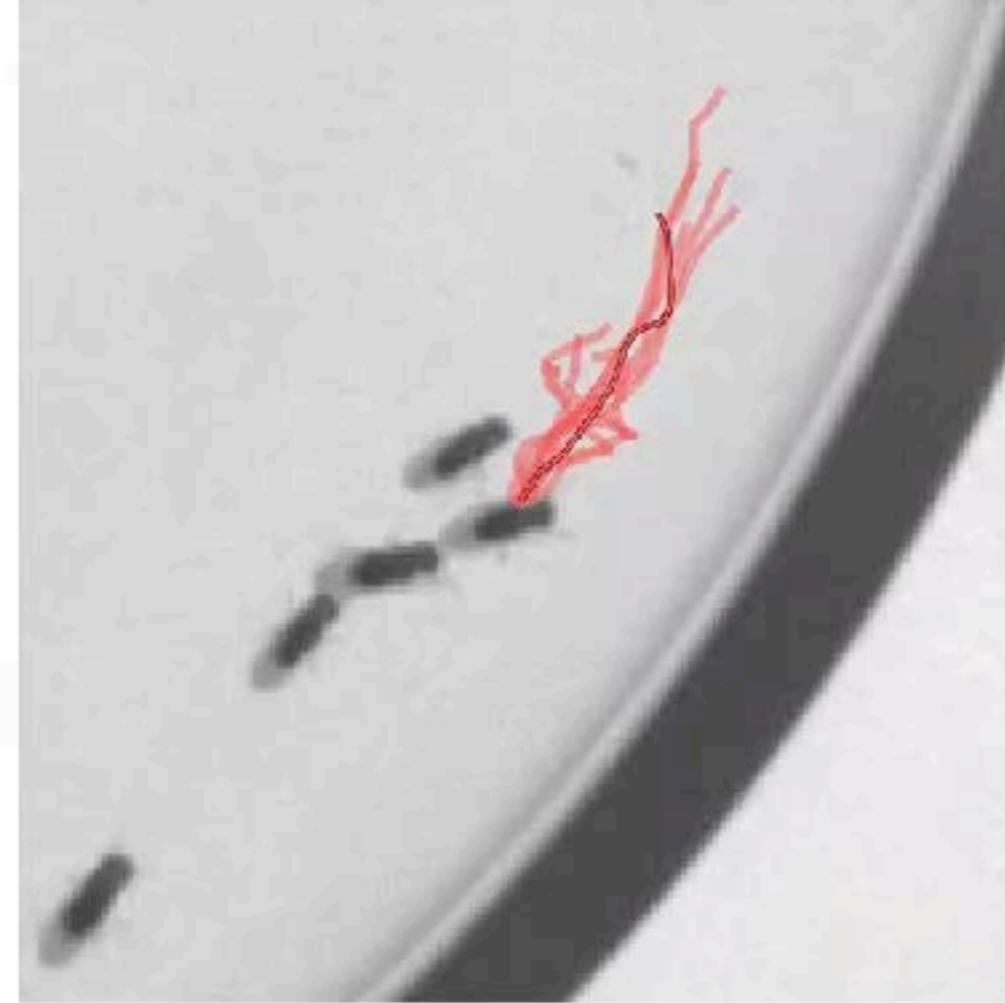
## tSNE dimensionality reduction

- female
- male
- right wing extension
- left wing extension

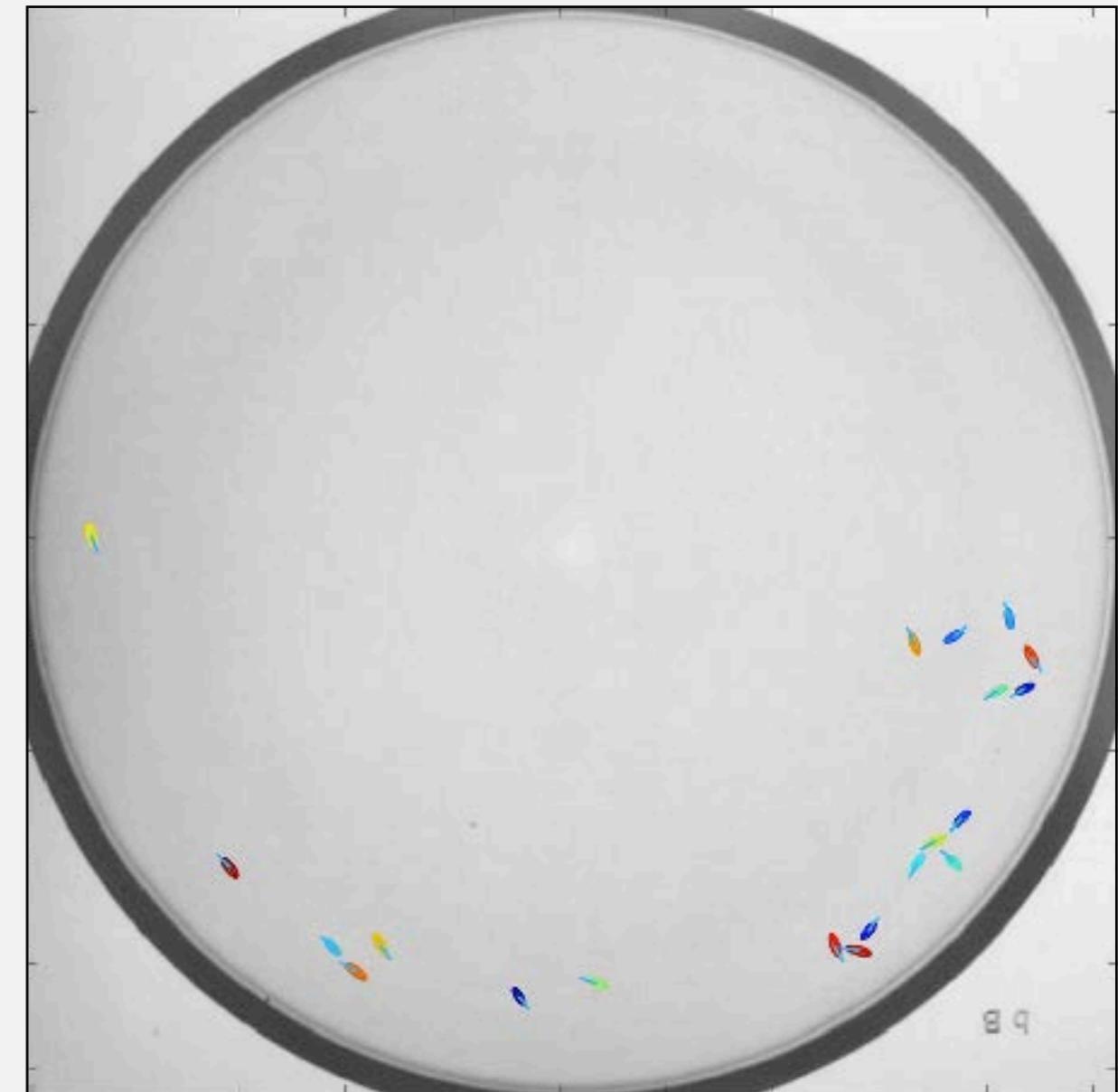
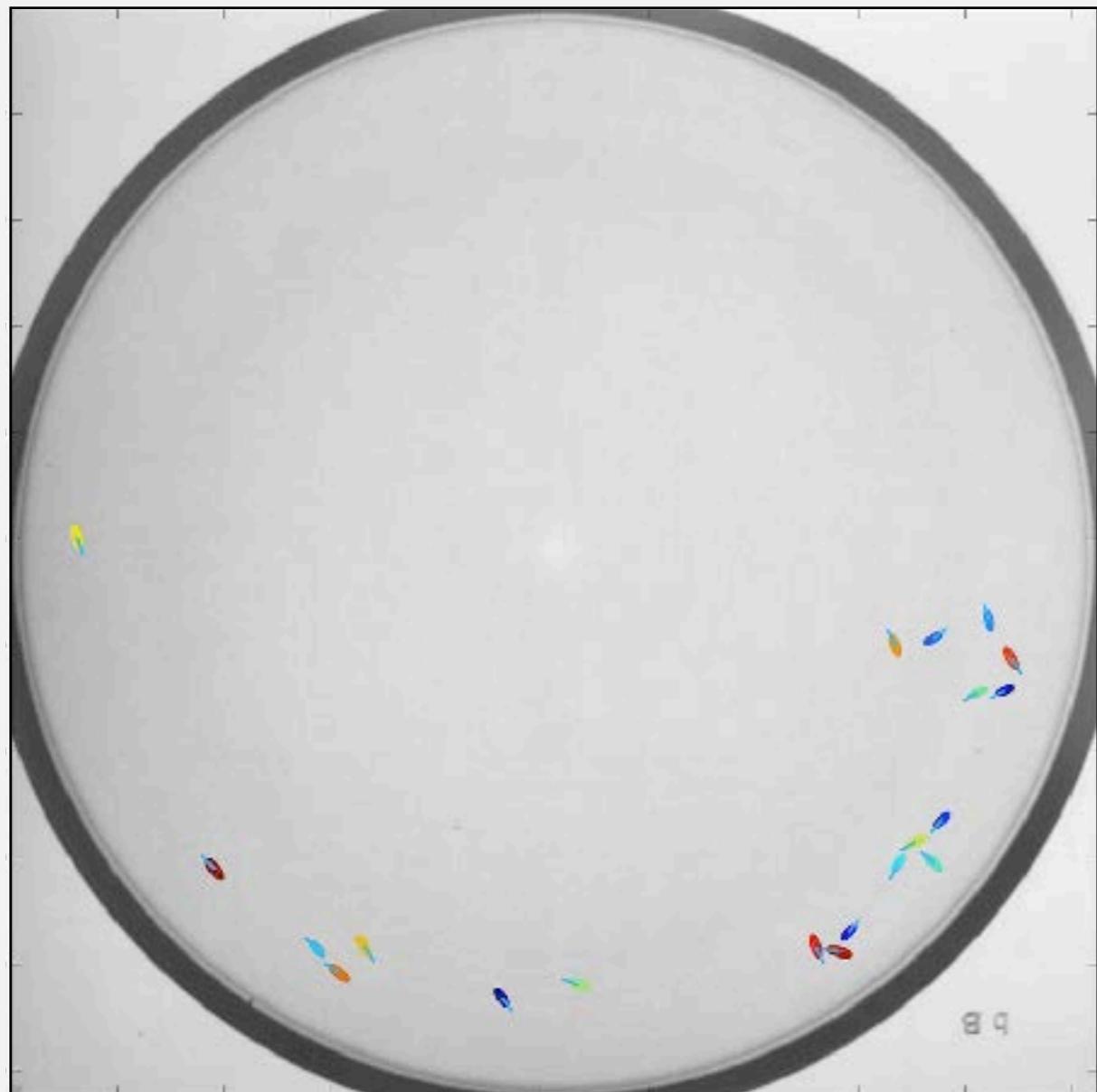
$h^l_i$ : hidden states of unsupervised model



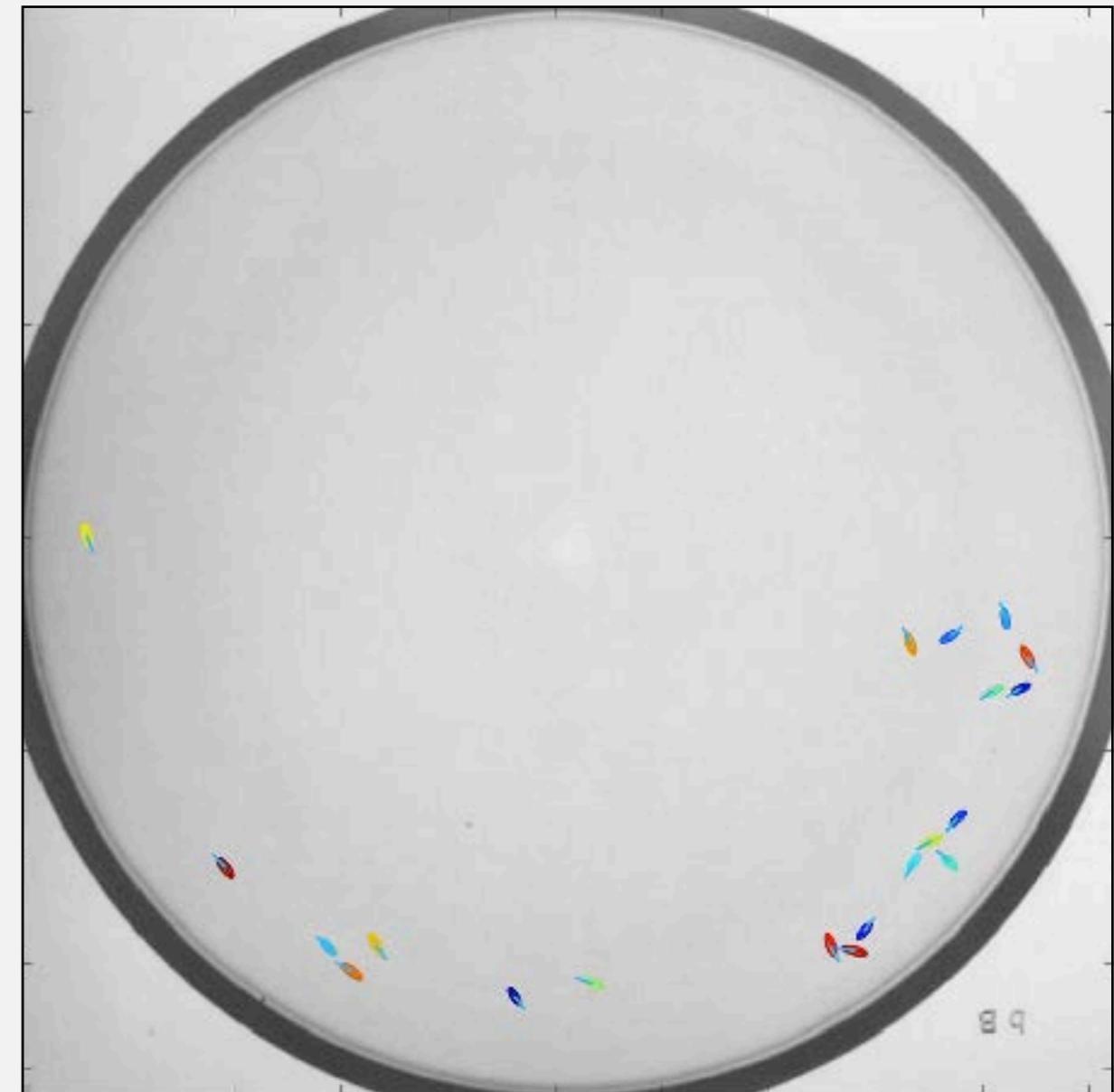
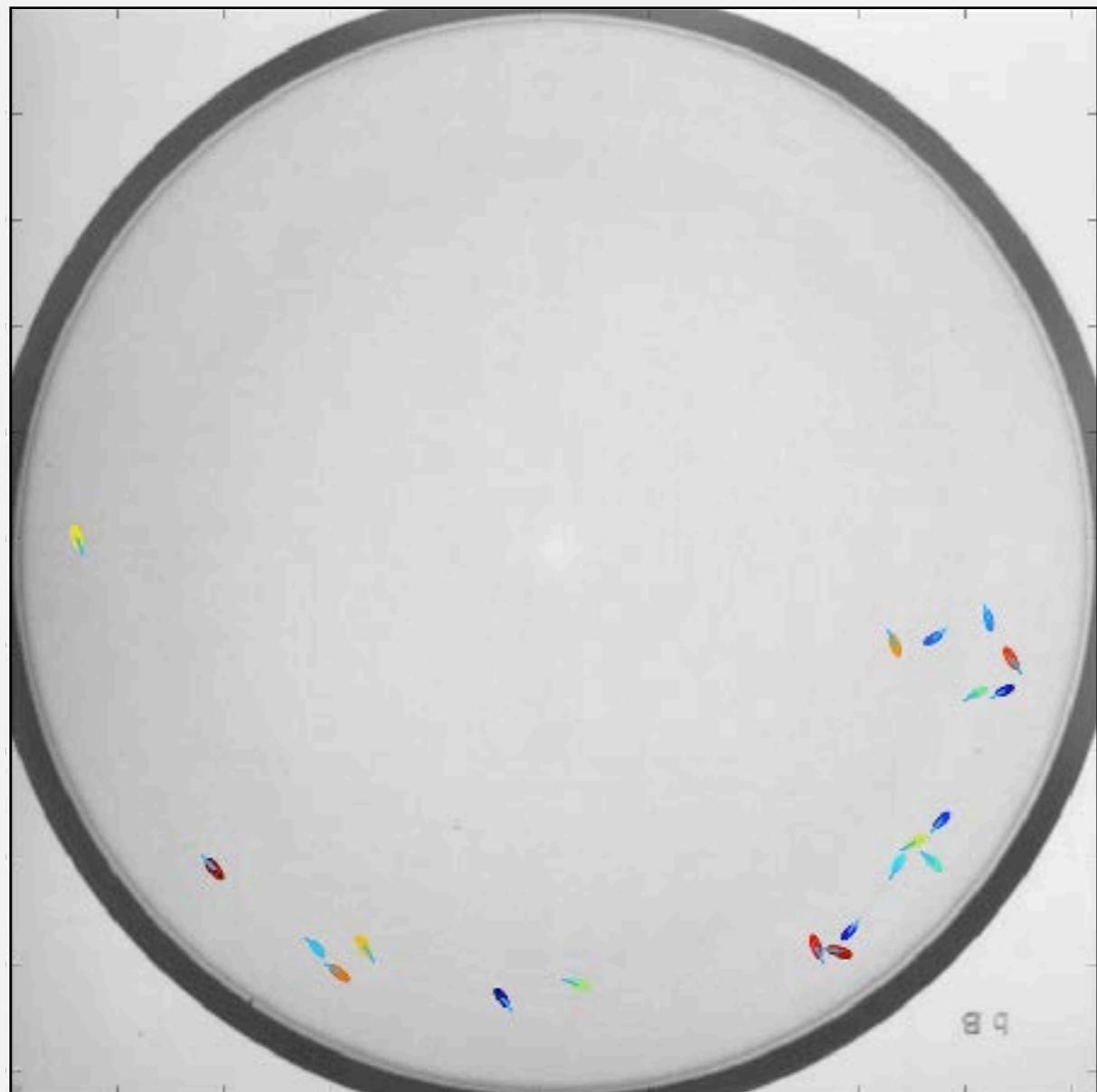




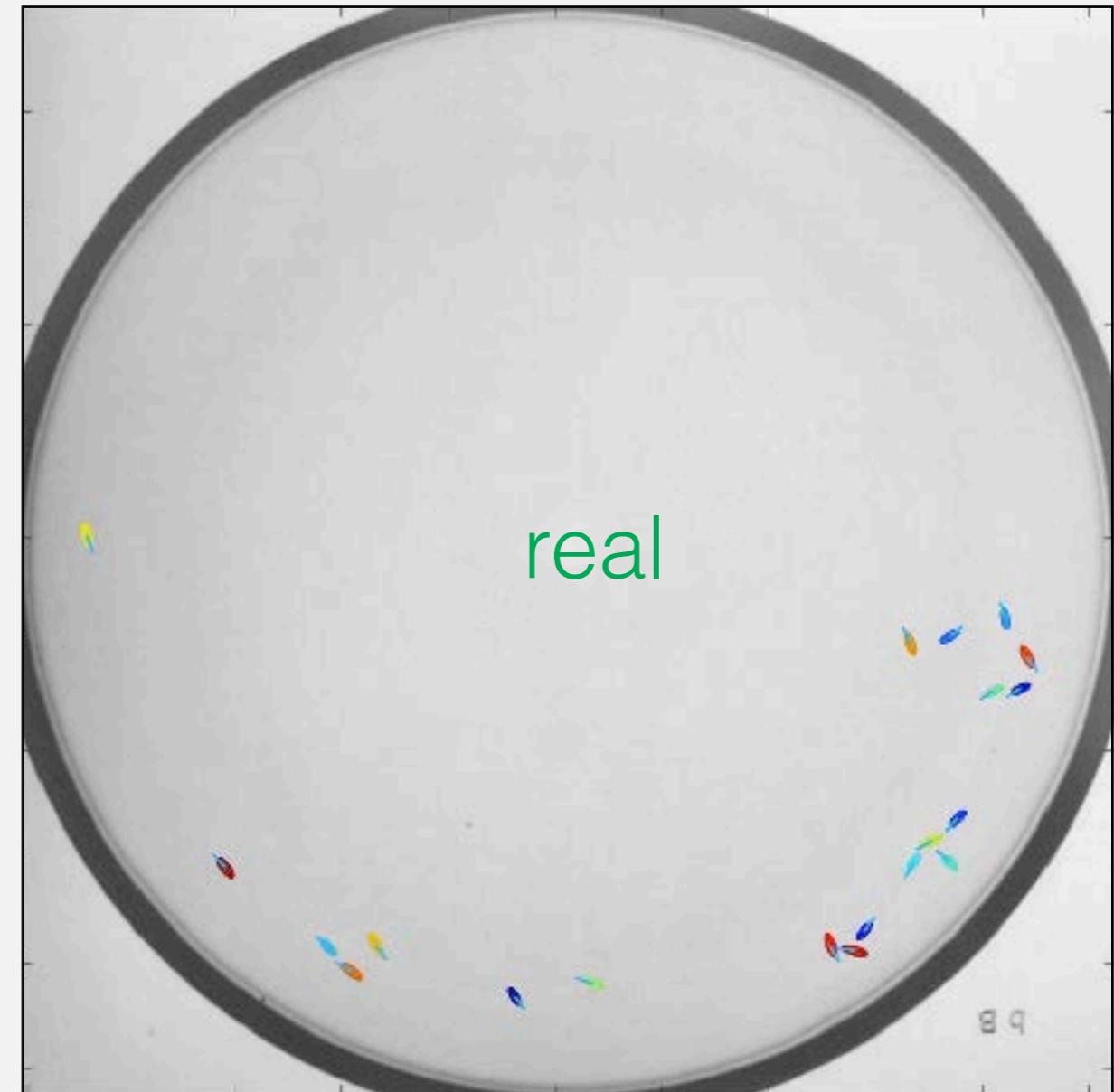
# Simulation: FlyBowl



# Simulation: FlyBowl



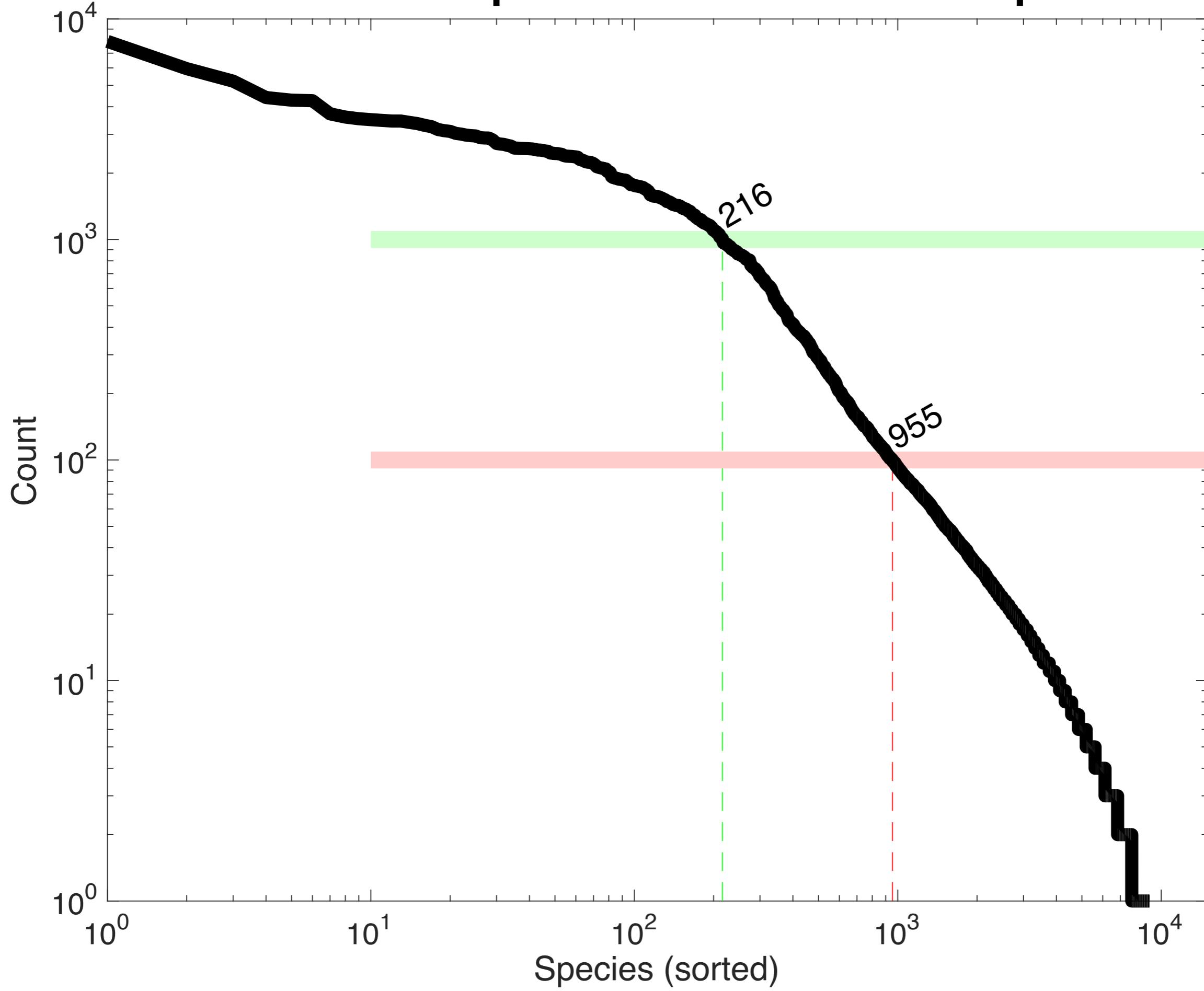
# Simulation: FlyBowl



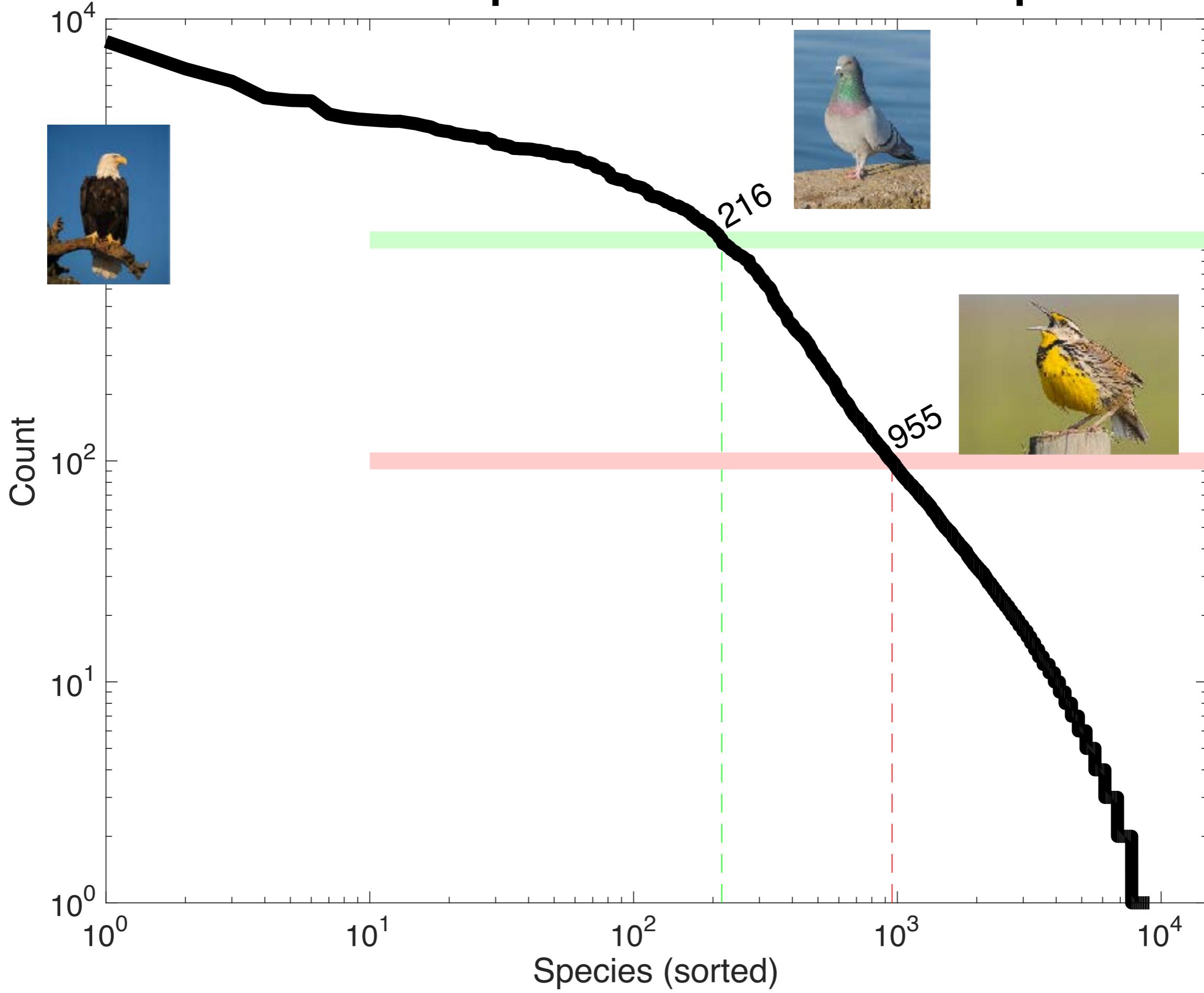
# Challenges

long tails ( $N > 0$ )

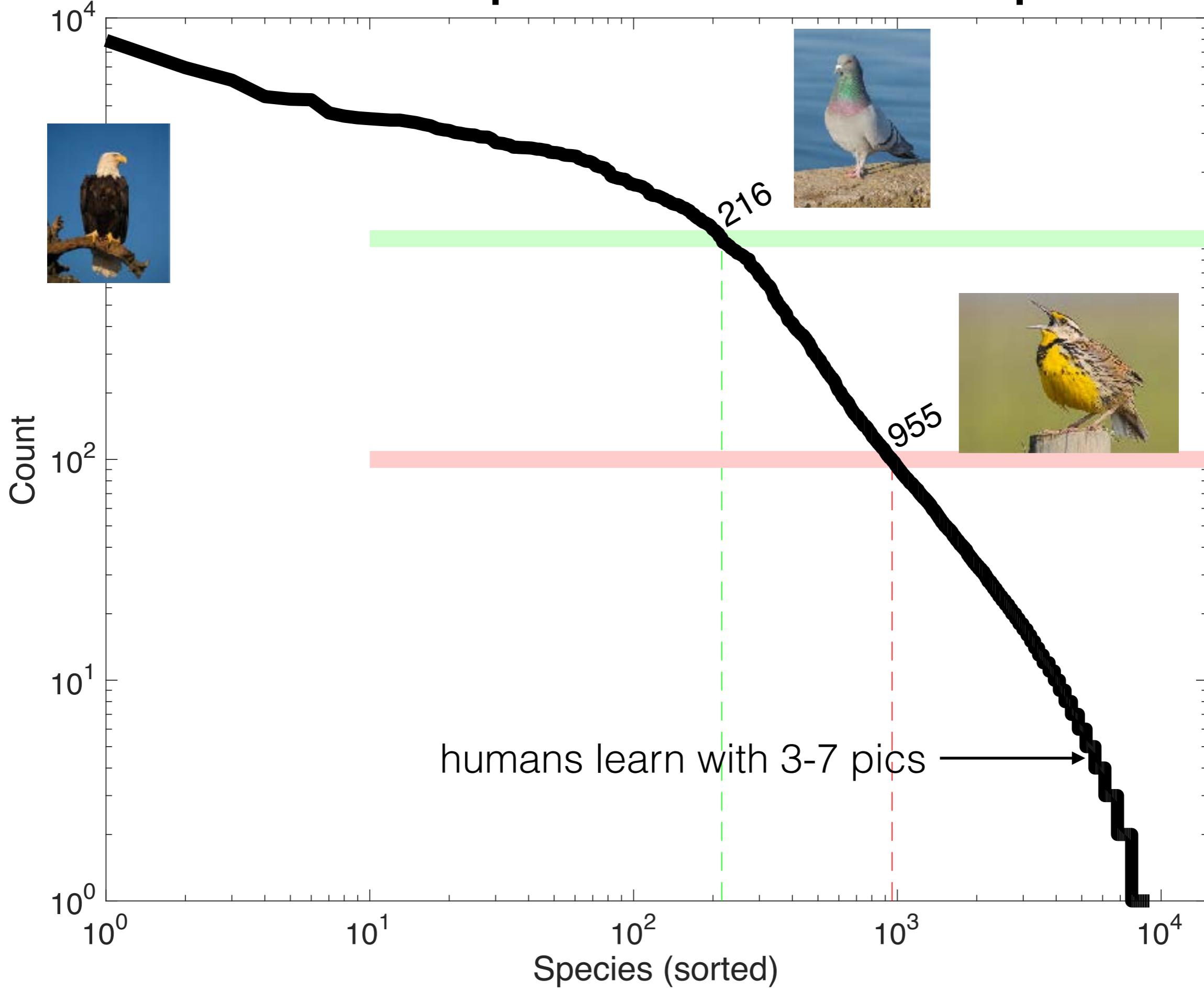
# eBird: 769561 specimens - 9031/15000 species



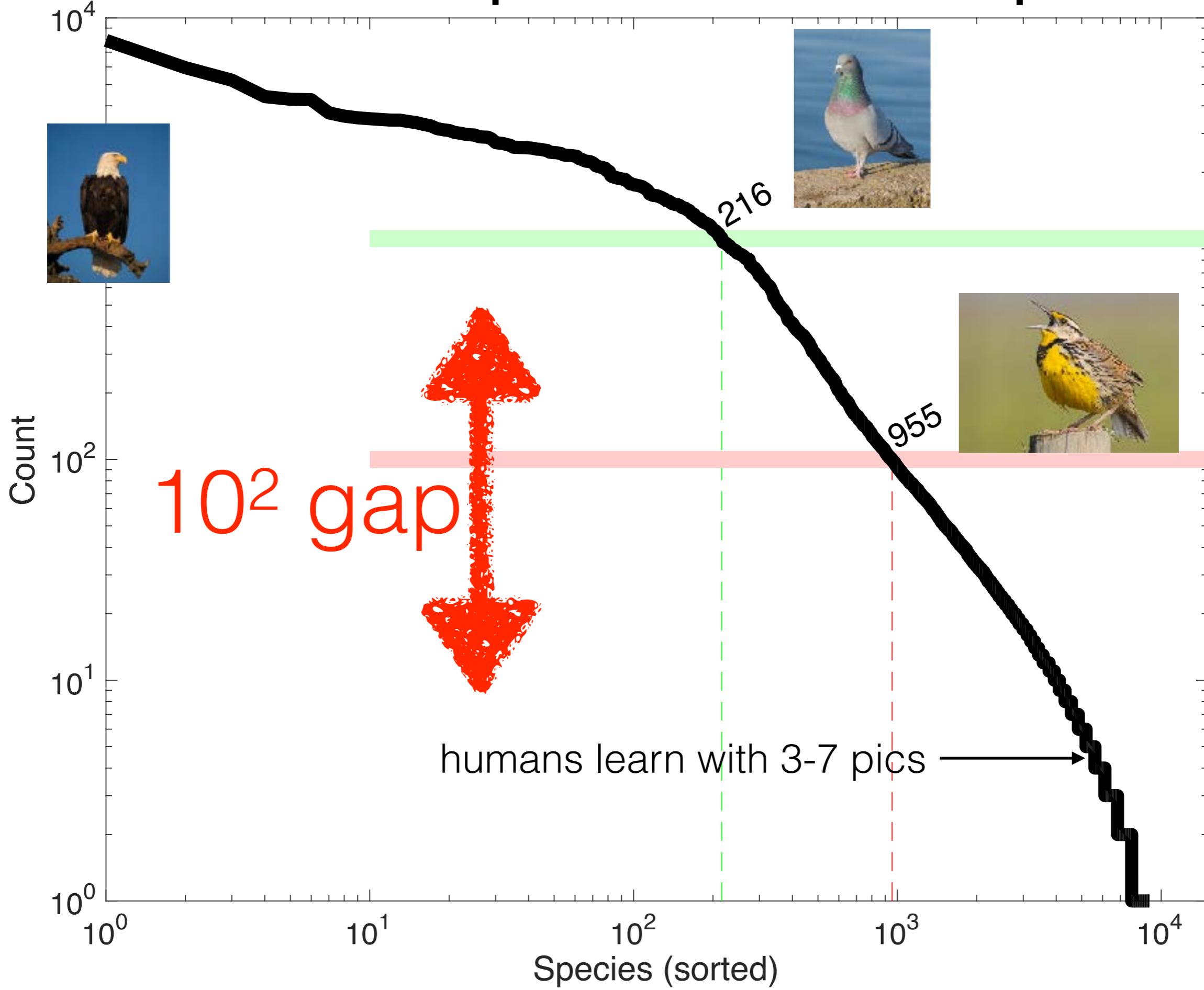
# eBird: 769561 specimens - 9031/15000 species



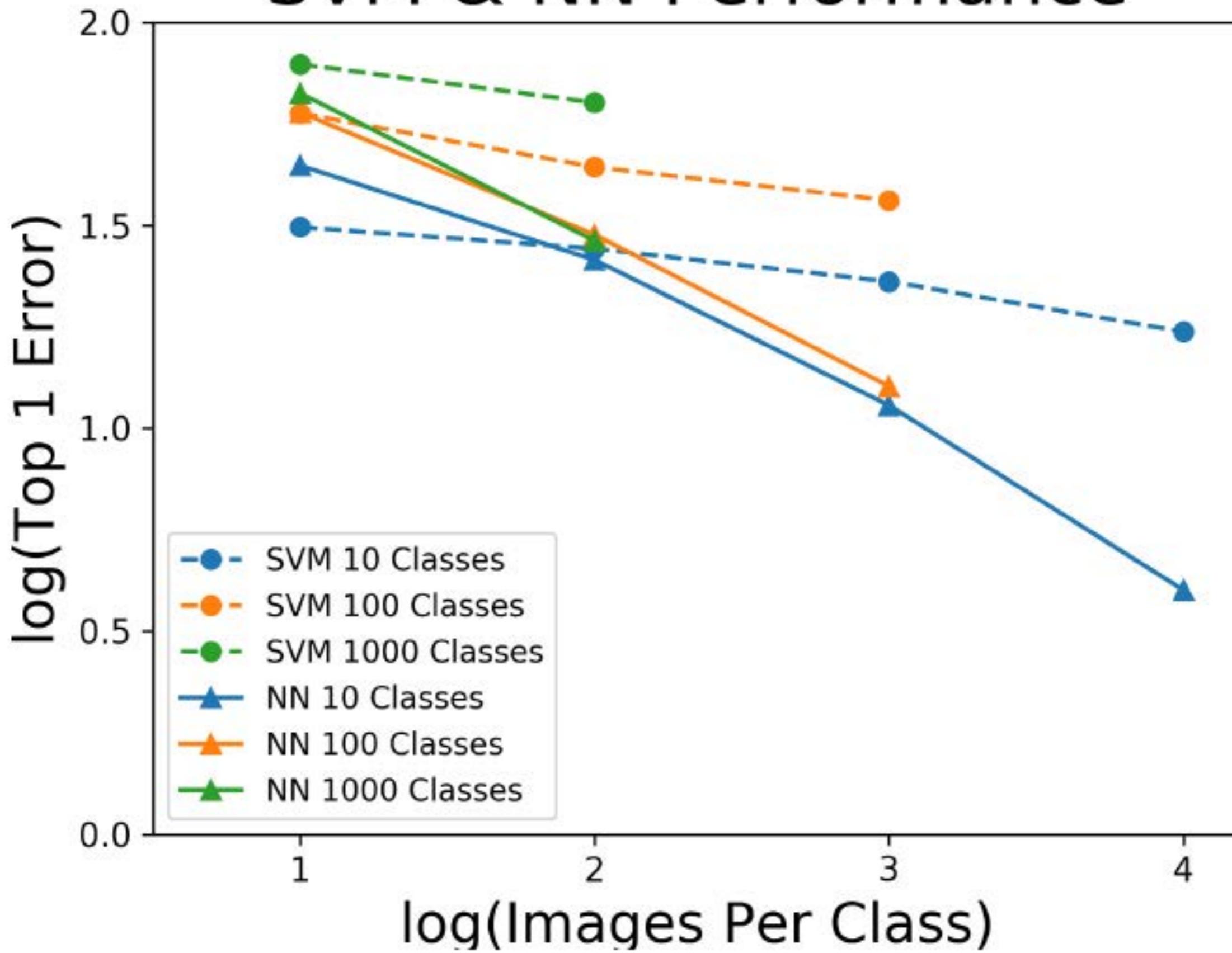
# eBird: 769561 specimens - 9031/15000 species



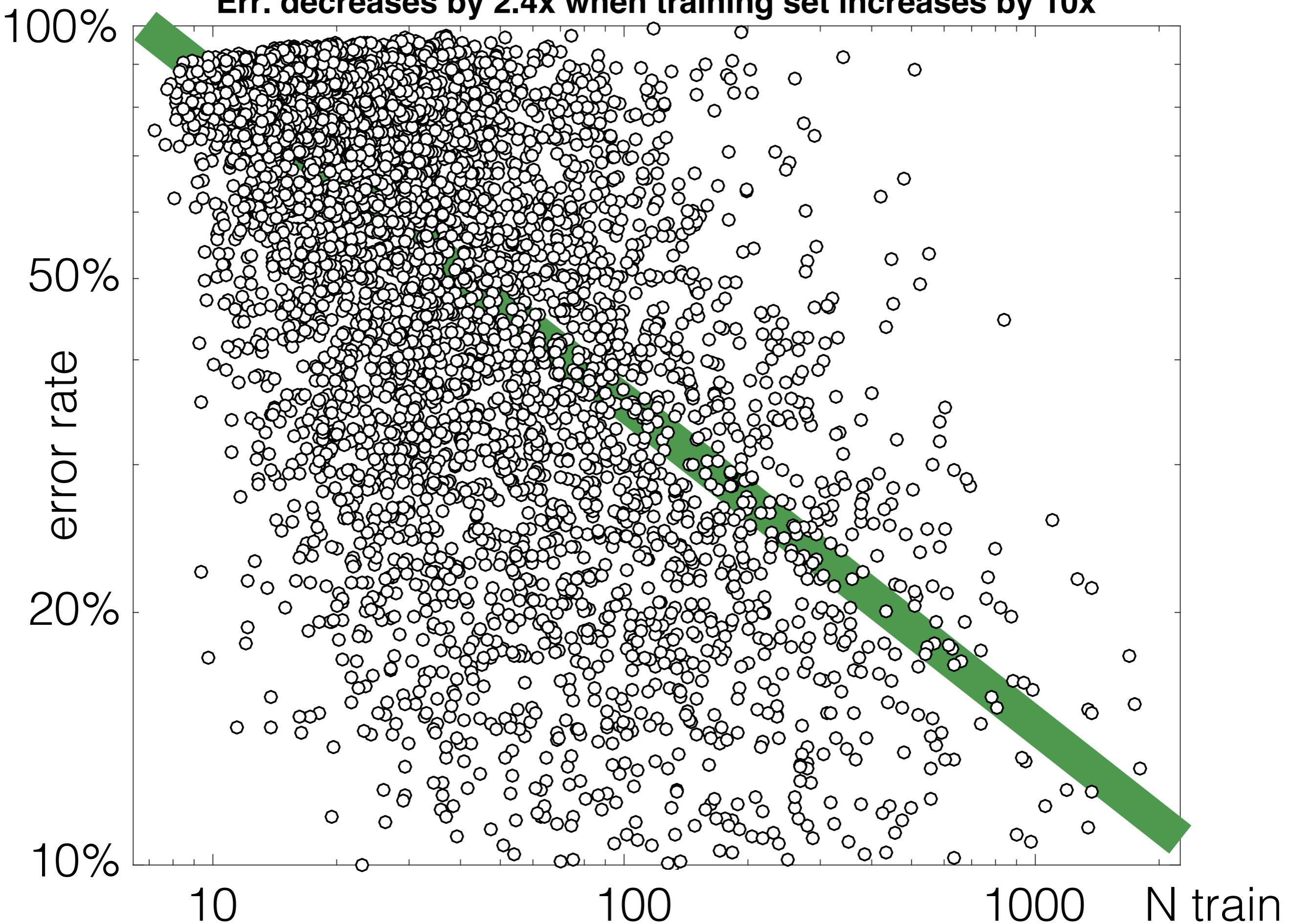
# eBird: 769561 specimens - 9031/15000 species



# SVM & NN Performance



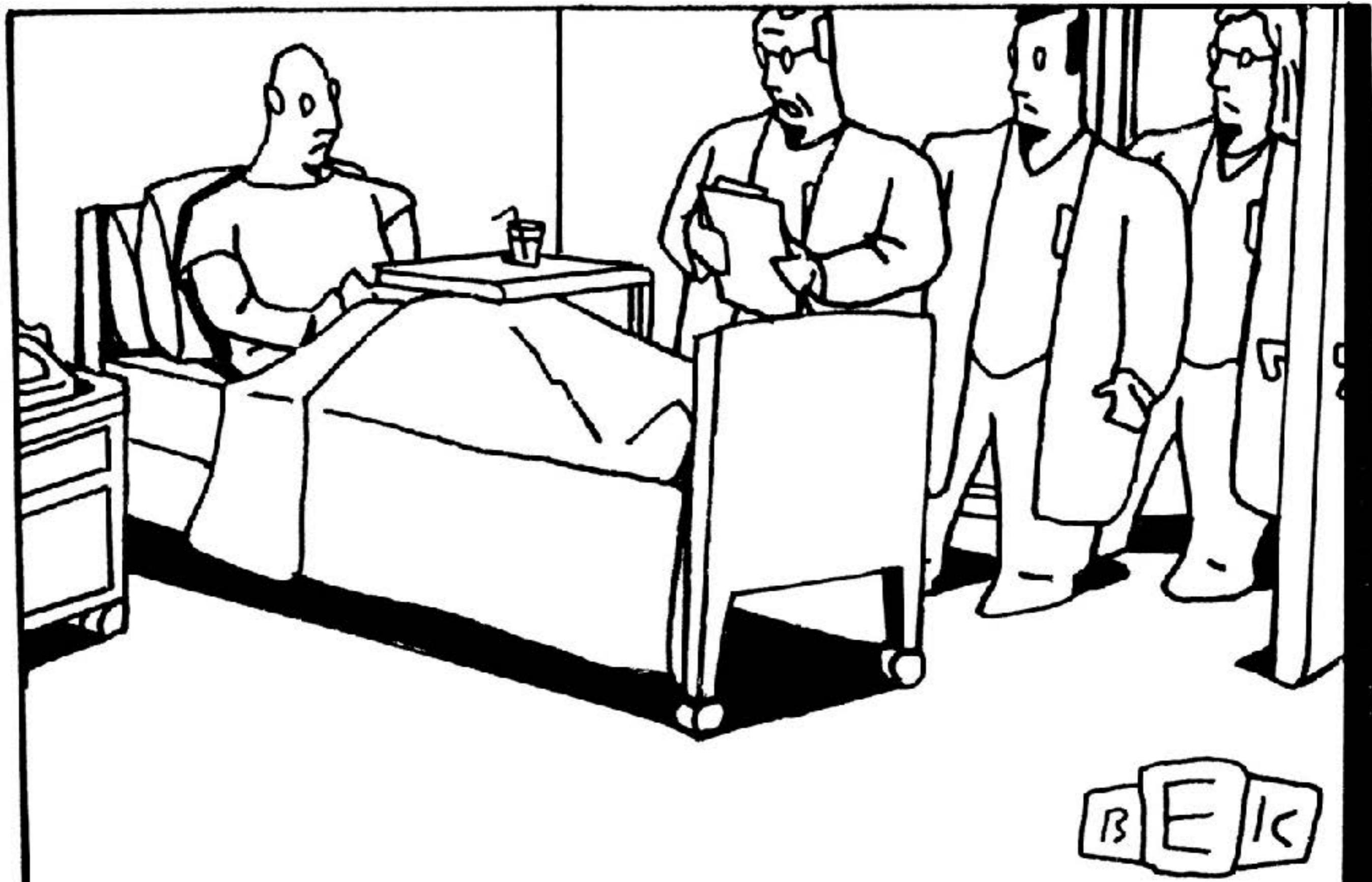
**Err. decreases by 2.4x when training set increases by 10x**



# Levels of understanding

- Memorization / recall
- Generalization / prediction
- Mechanisms / intervention

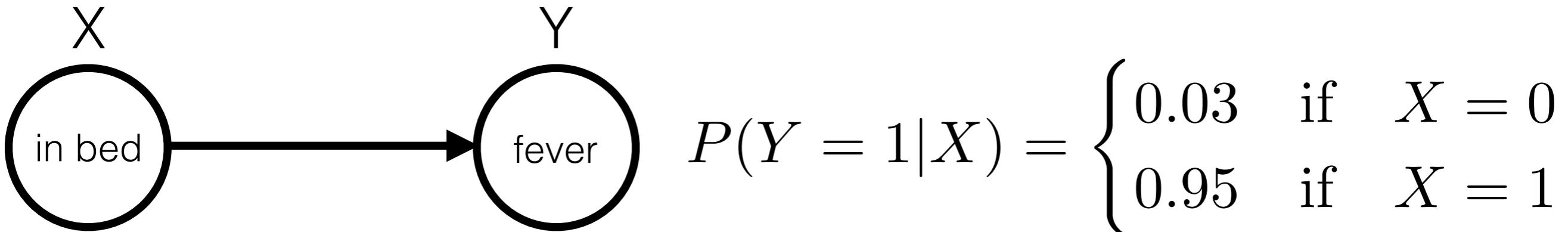
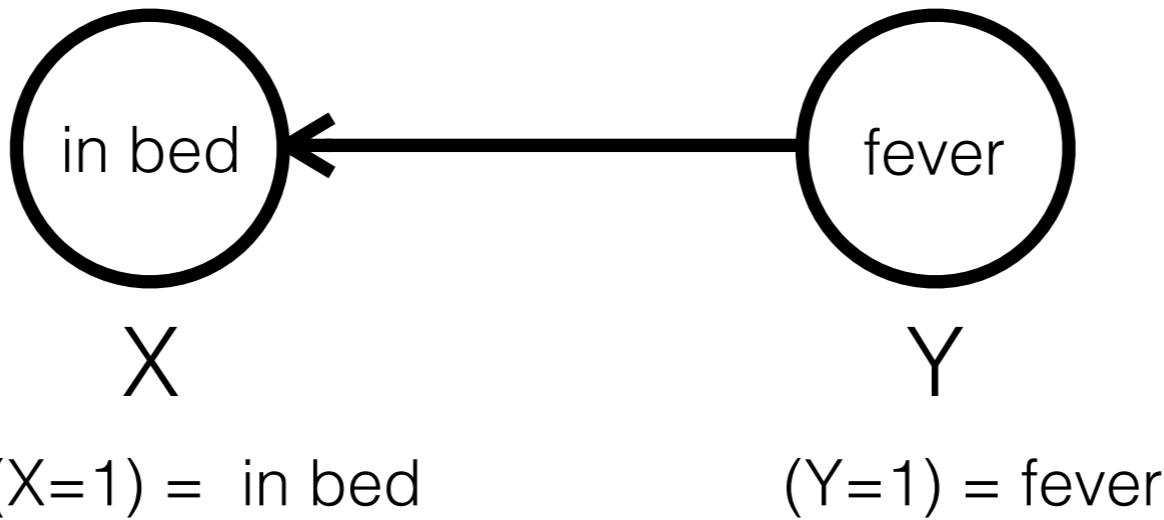
# Correlation vs causation



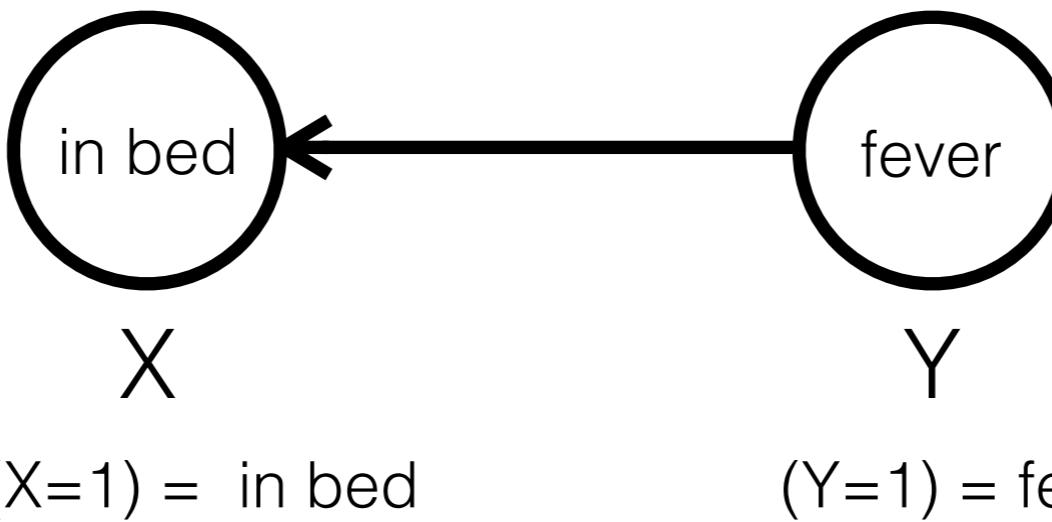


$$P(Y = 1|X) = \begin{cases} 0.03 & \text{if } X = 0 \\ 0.95 & \text{if } X = 1 \end{cases}$$

# Causation



# Causation



cause  
←  
dependence  
←

# Correlation

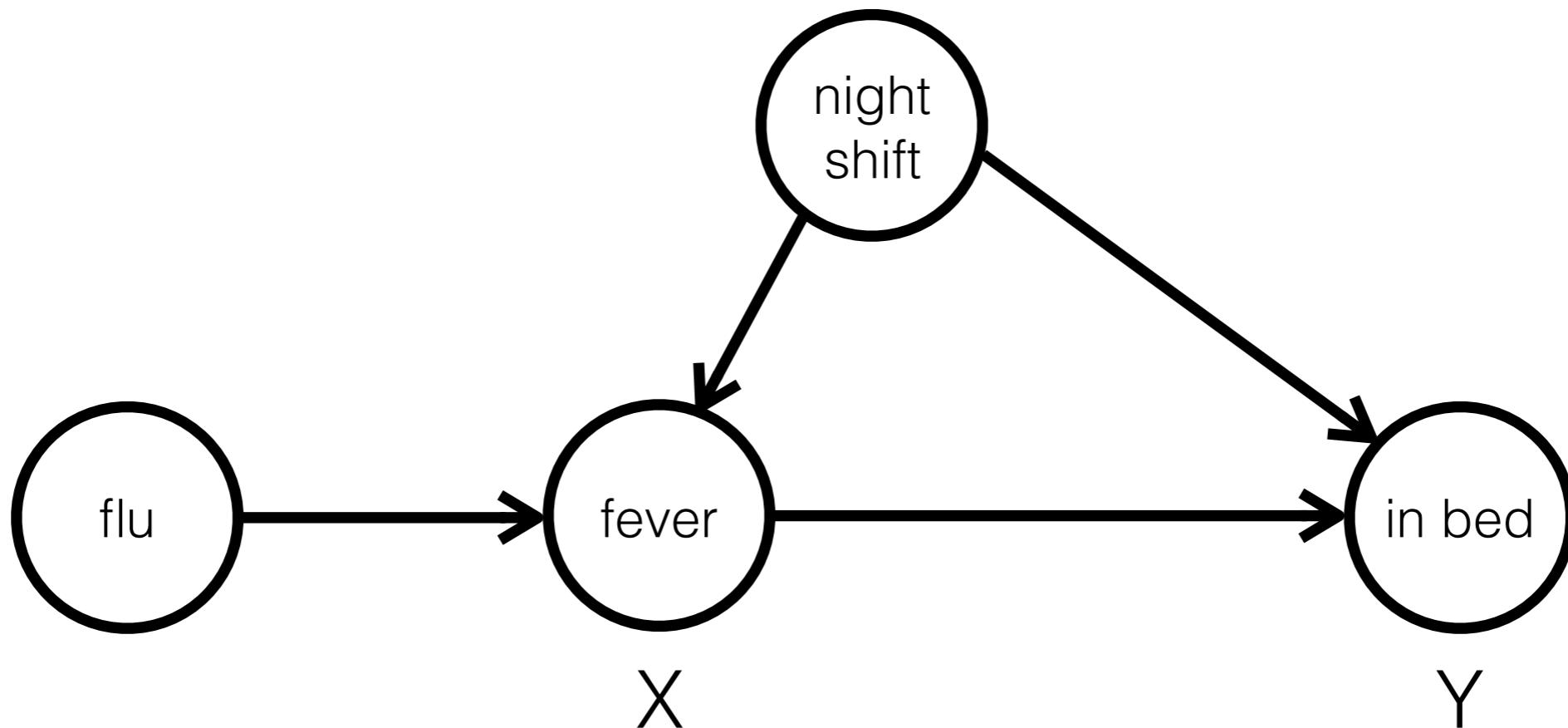


$$P(X = 1|Y) = \begin{cases} 0.01 & \text{if } Y = 0 \\ 0.9 & \text{if } Y = 1 \end{cases}$$



$$P(Y = 1|X) = \begin{cases} 0.03 & \text{if } X = 0 \\ 0.95 & \text{if } X = 1 \end{cases}$$

# Definition of causation

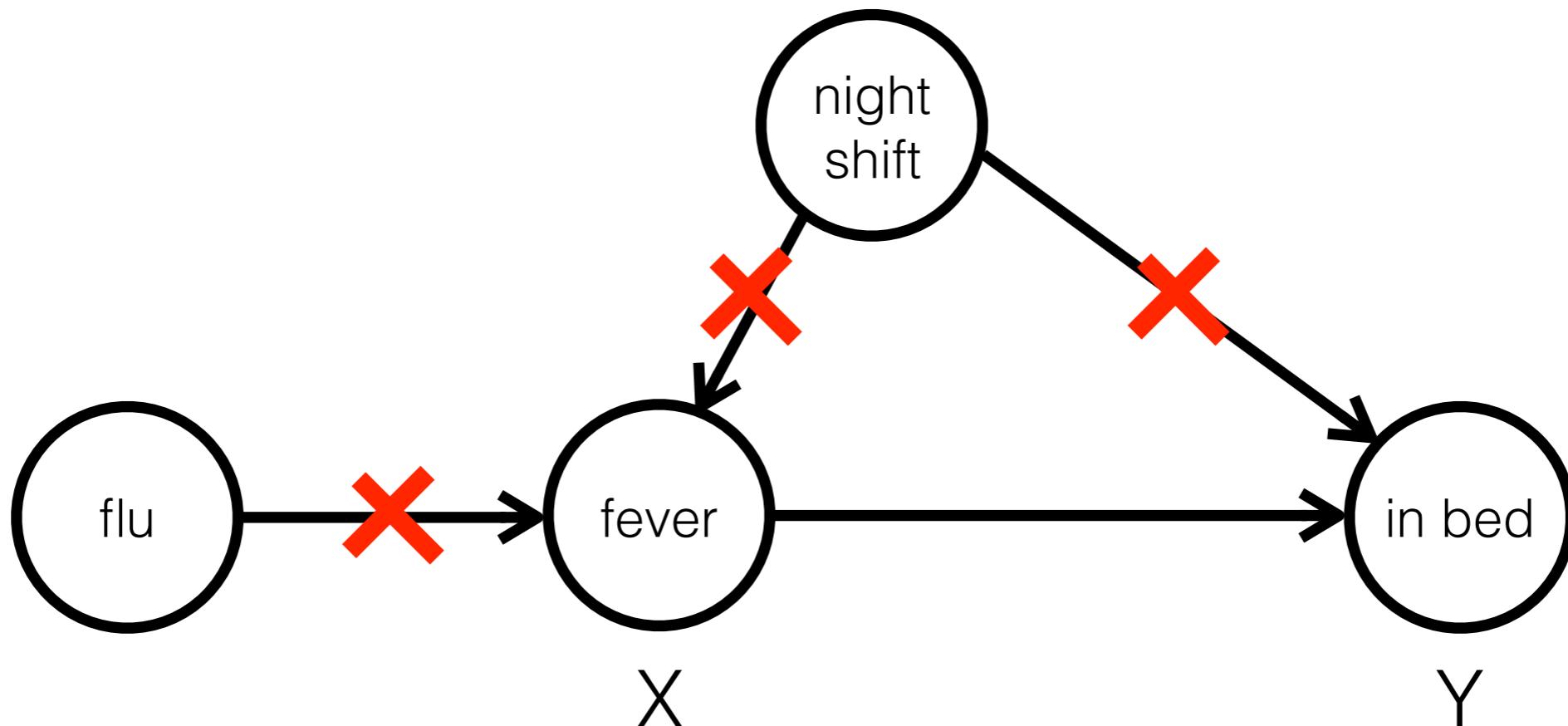


$$X \not\perp\!\!\!\perp Y \Leftrightarrow P(Y | X) \neq P(Y)$$

$$X \rightarrow Y \Leftrightarrow P(Y | \text{do}(X)) \neq P(Y)$$

[Pearl 2000]

# Definition of causation

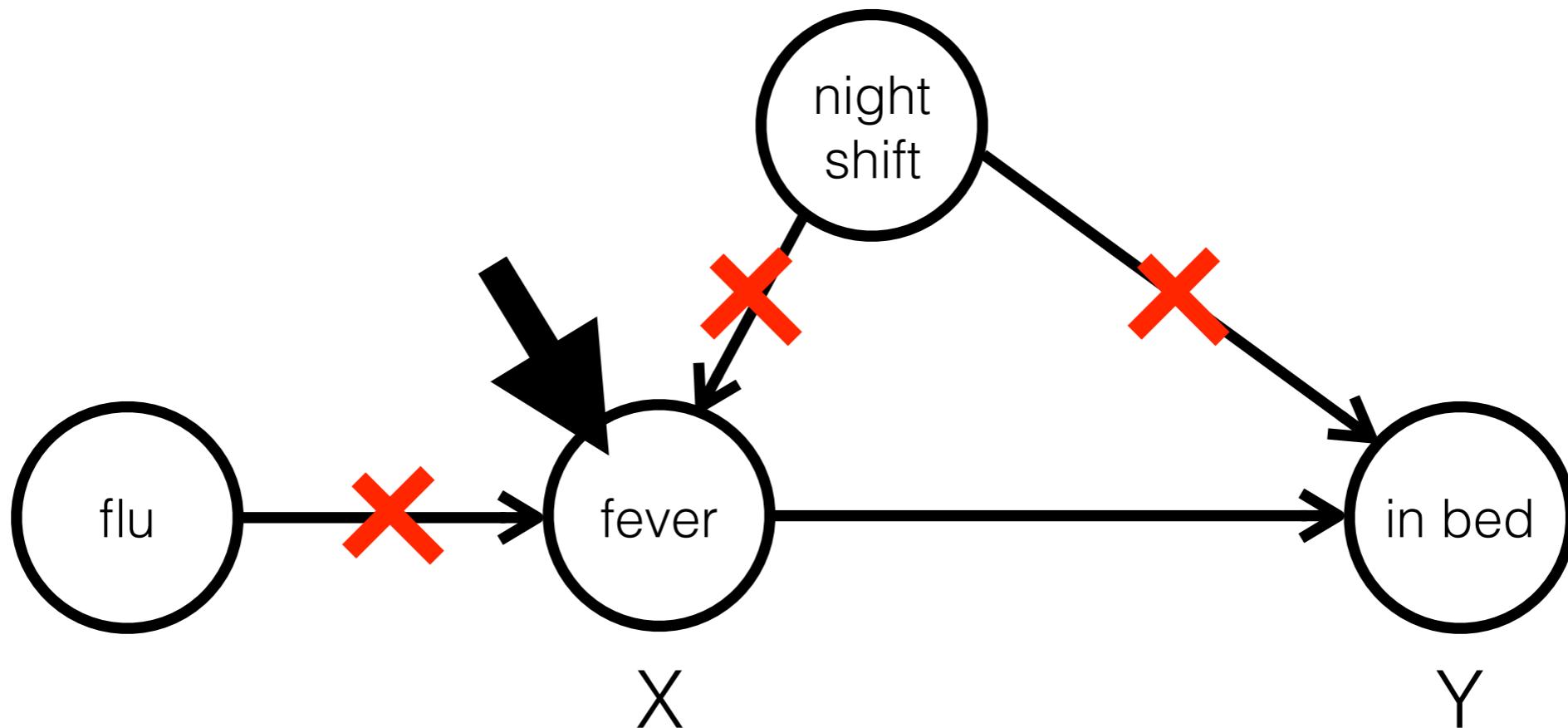


$$X \not\perp\!\!\!\perp Y \Leftrightarrow P(Y | X) \neq P(Y)$$

$$X \rightarrow Y \Leftrightarrow P(Y | \text{do}(X)) \neq P(Y)$$

[Pearl 2000]

# Definition of causation

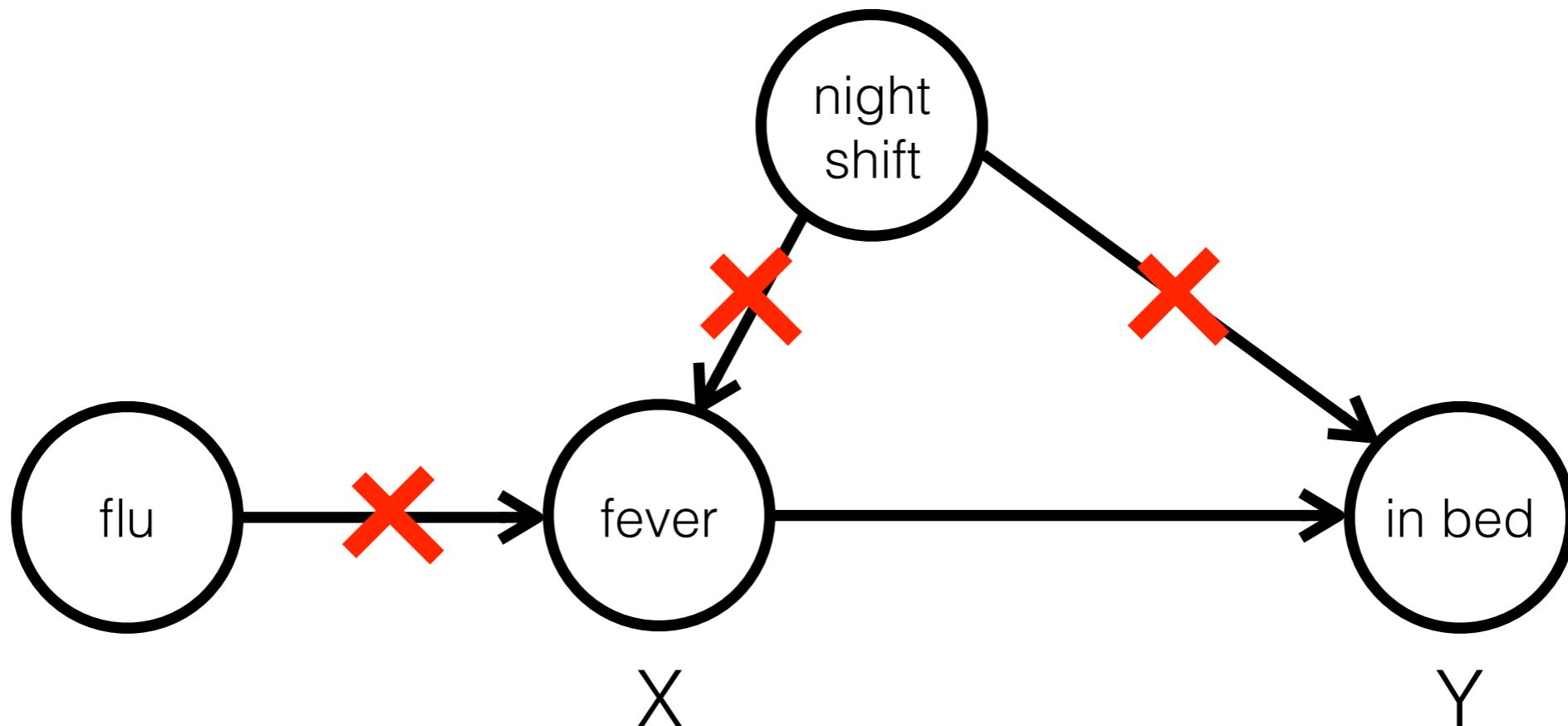


$$X \not\perp\!\!\!\perp Y \Leftrightarrow P(Y | X) \neq P(Y)$$

$$X \rightarrow Y \Leftrightarrow P(Y | \text{do}(X)) \neq P(Y)$$

[Pearl 2000]

# Definition of causation

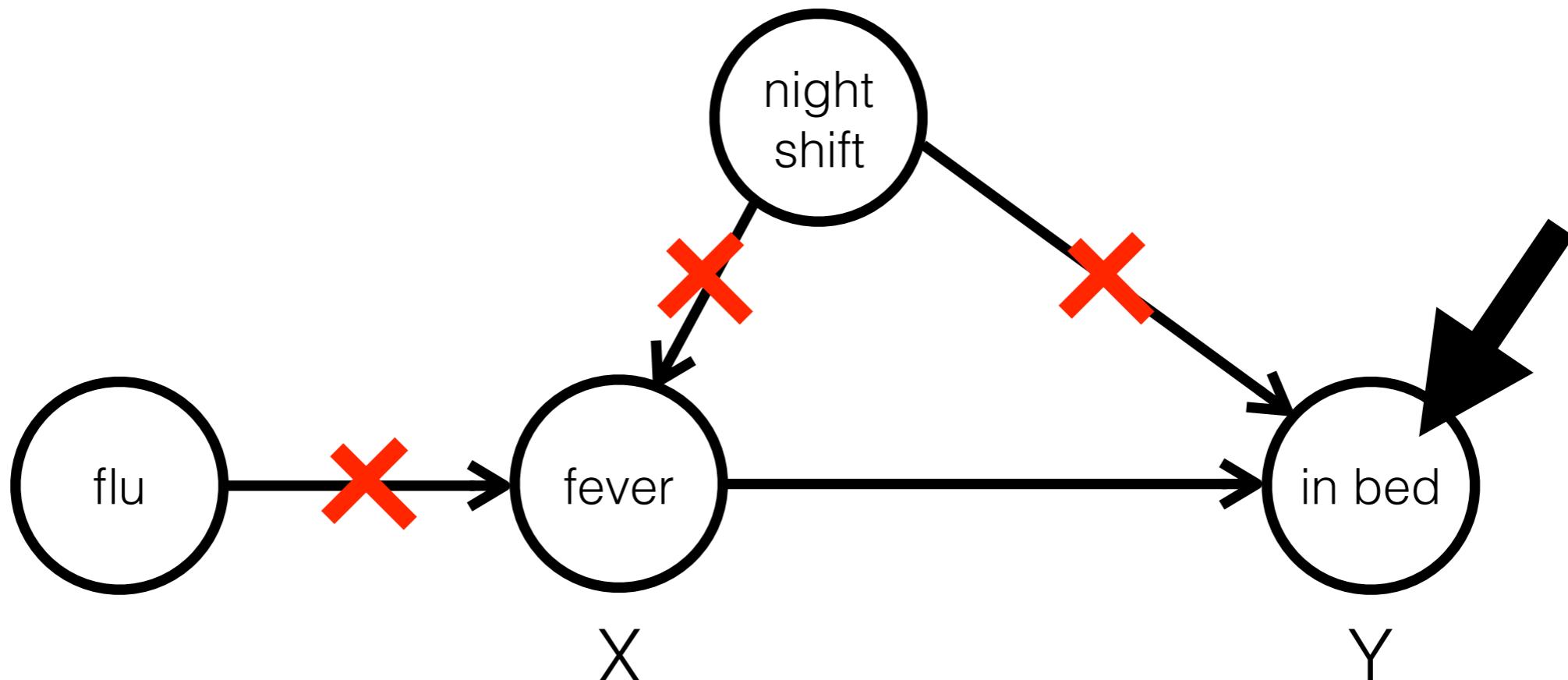


$$X \not\perp\!\!\!\perp Y \Leftrightarrow P(Y | X) \neq P(Y)$$

$$X \rightarrow Y \Leftrightarrow P(Y | \text{do}(X)) \neq P(Y)$$

[Pearl 2000]

# Definition of causation

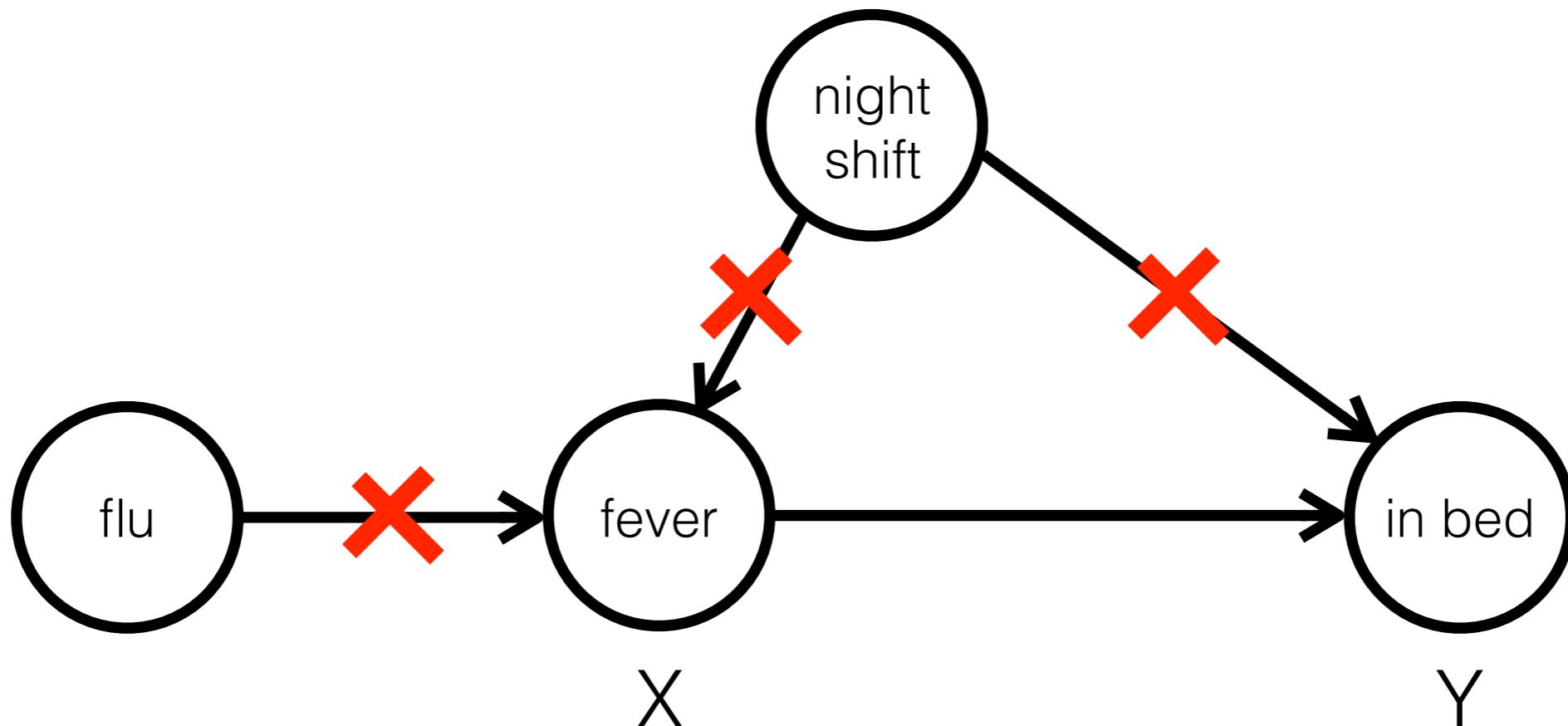


$$X \not\perp\!\!\!\perp Y \Leftrightarrow P(Y | X) \neq P(Y)$$

$$X \rightarrow Y \Leftrightarrow P(Y | \text{do}(X)) \neq P(Y)$$

[Pearl 2000]

# Definition of causation

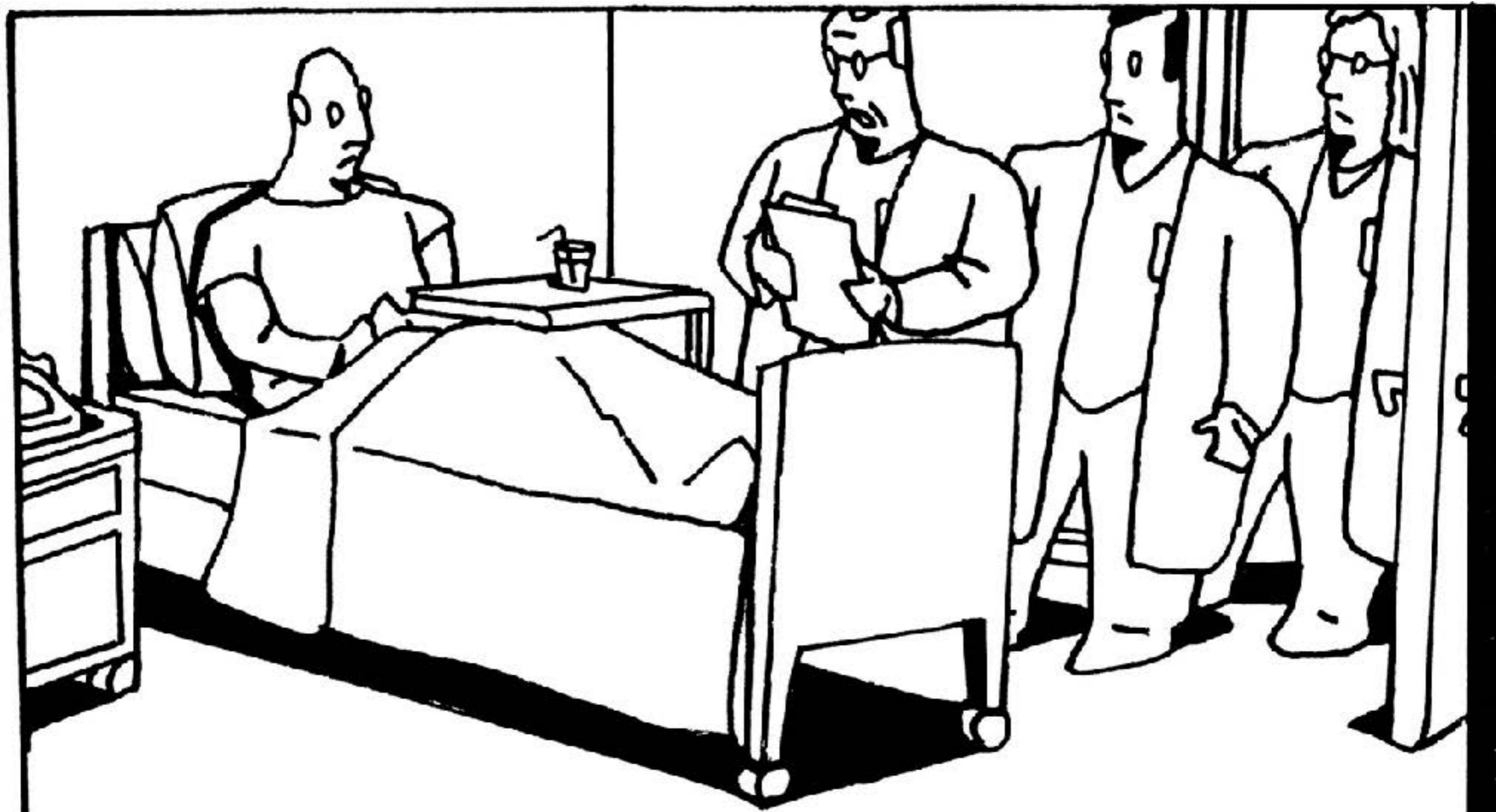


$$X \not\perp\!\!\!\perp Y \Leftrightarrow P(Y | X) \neq P(Y)$$

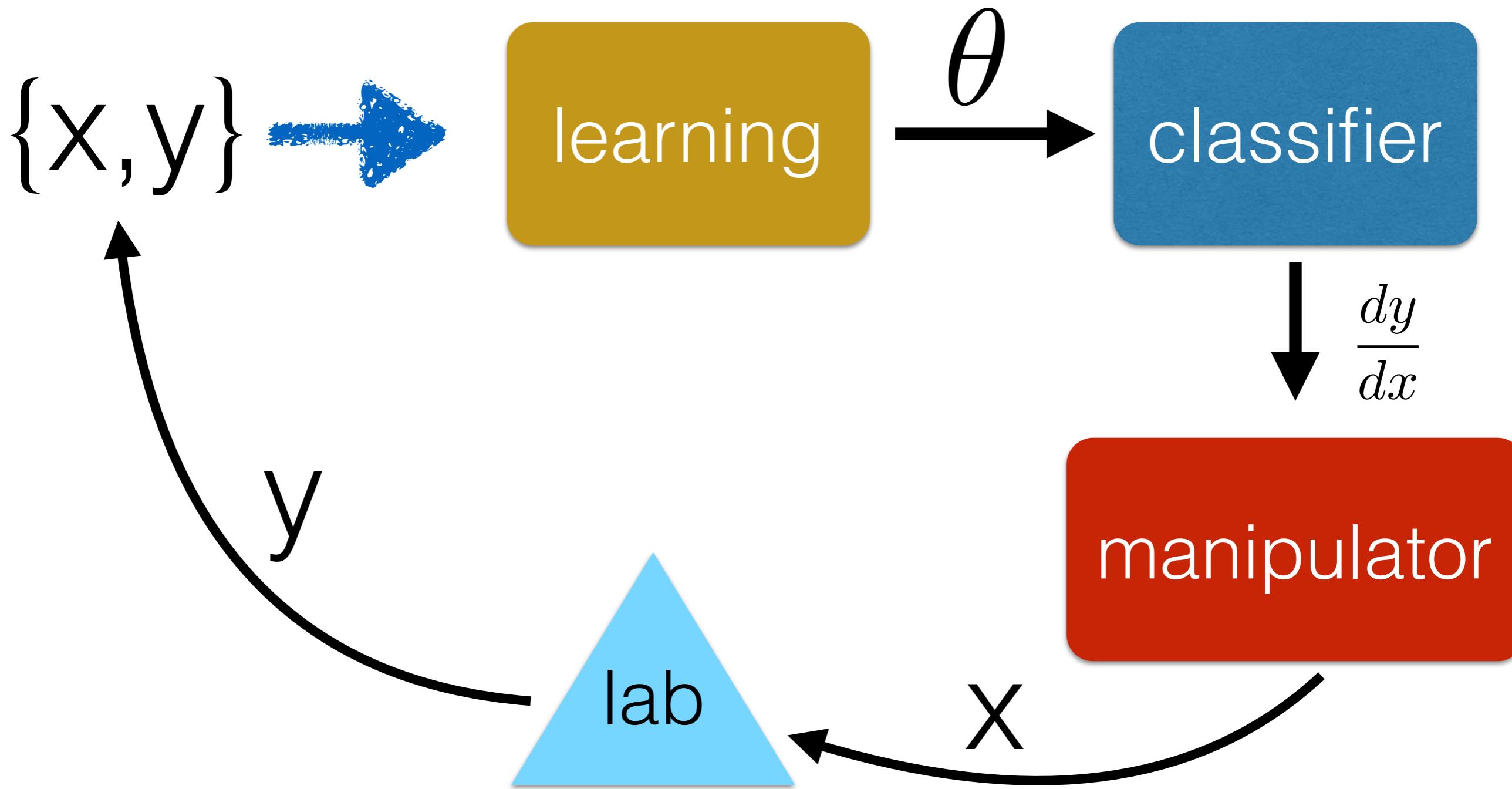
$$X \rightarrow Y \Leftrightarrow P(Y | \text{do}(X)) \neq P(Y)$$

[Pearl 2000]

# Prediction vs intervention

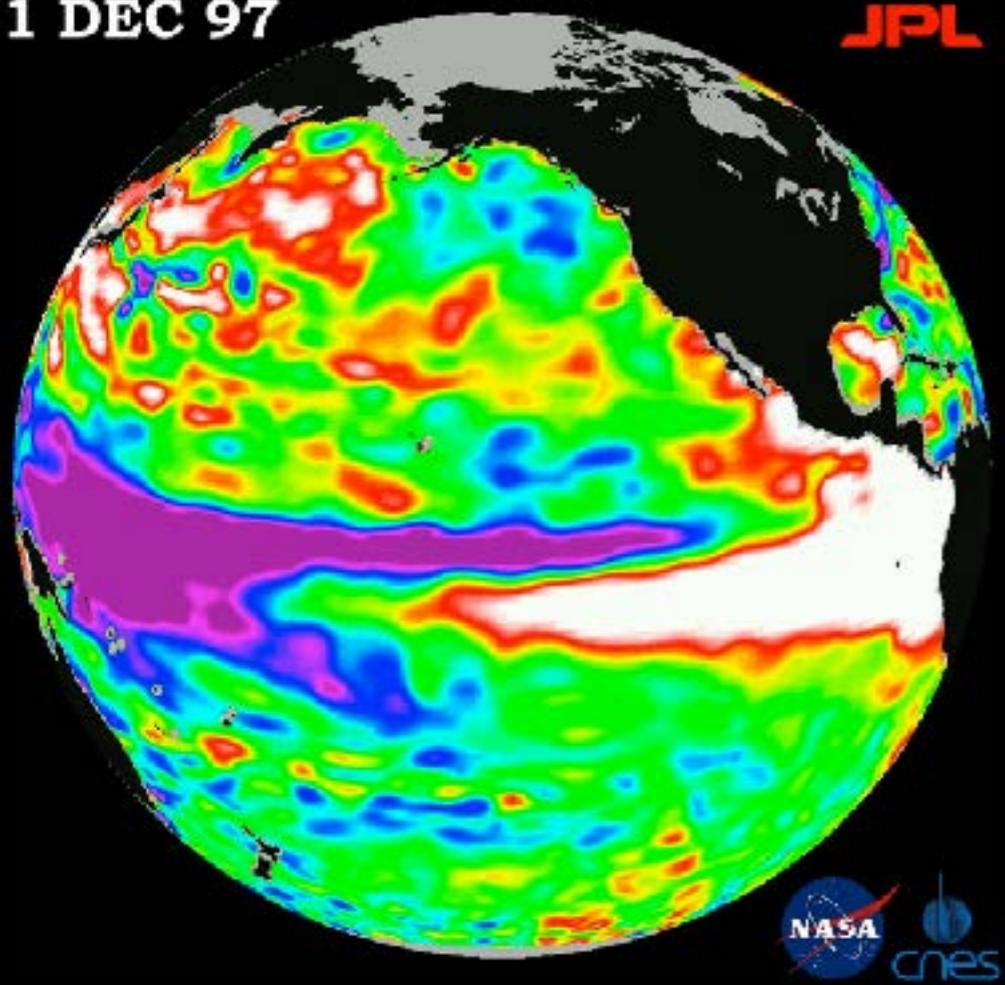


B E IC

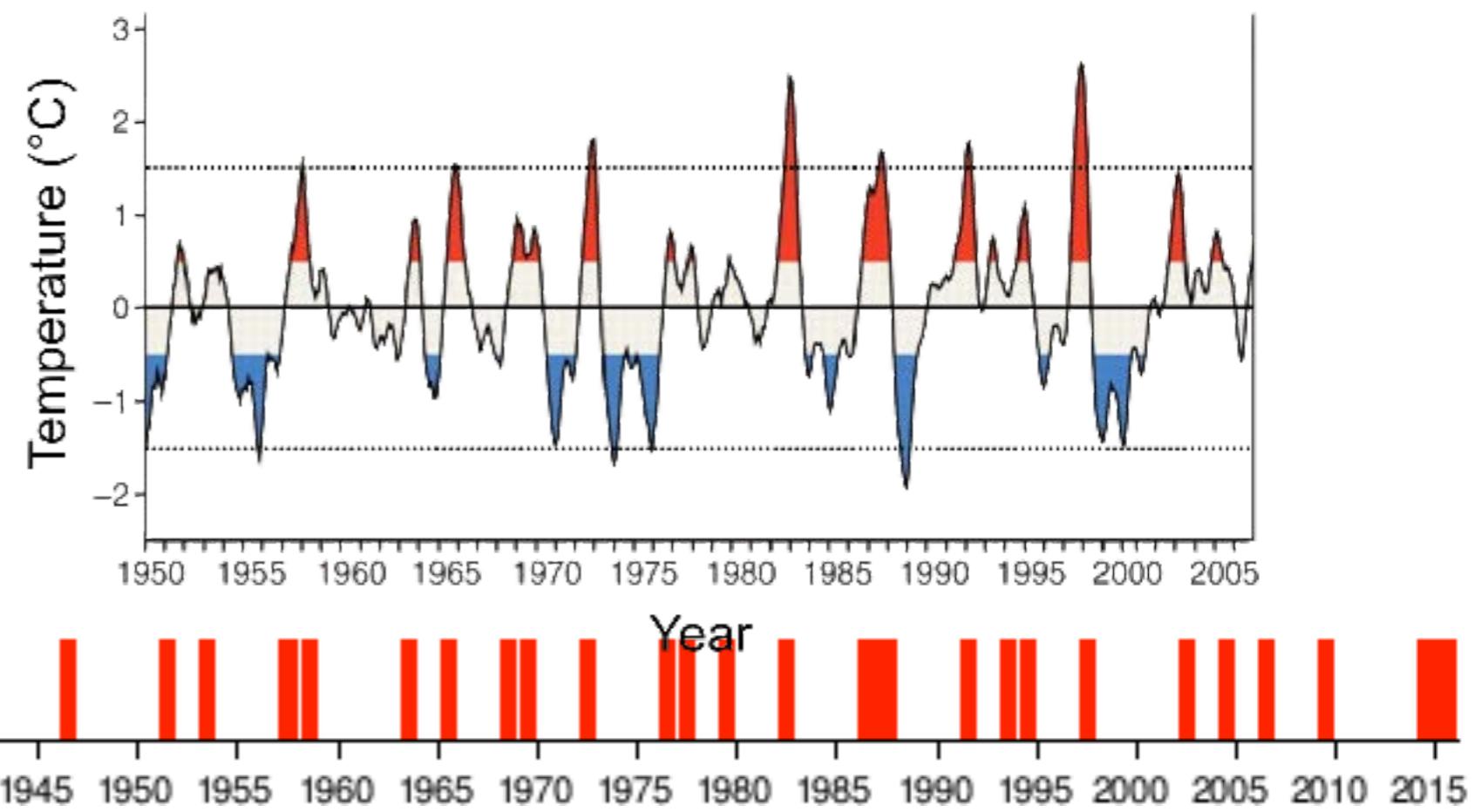
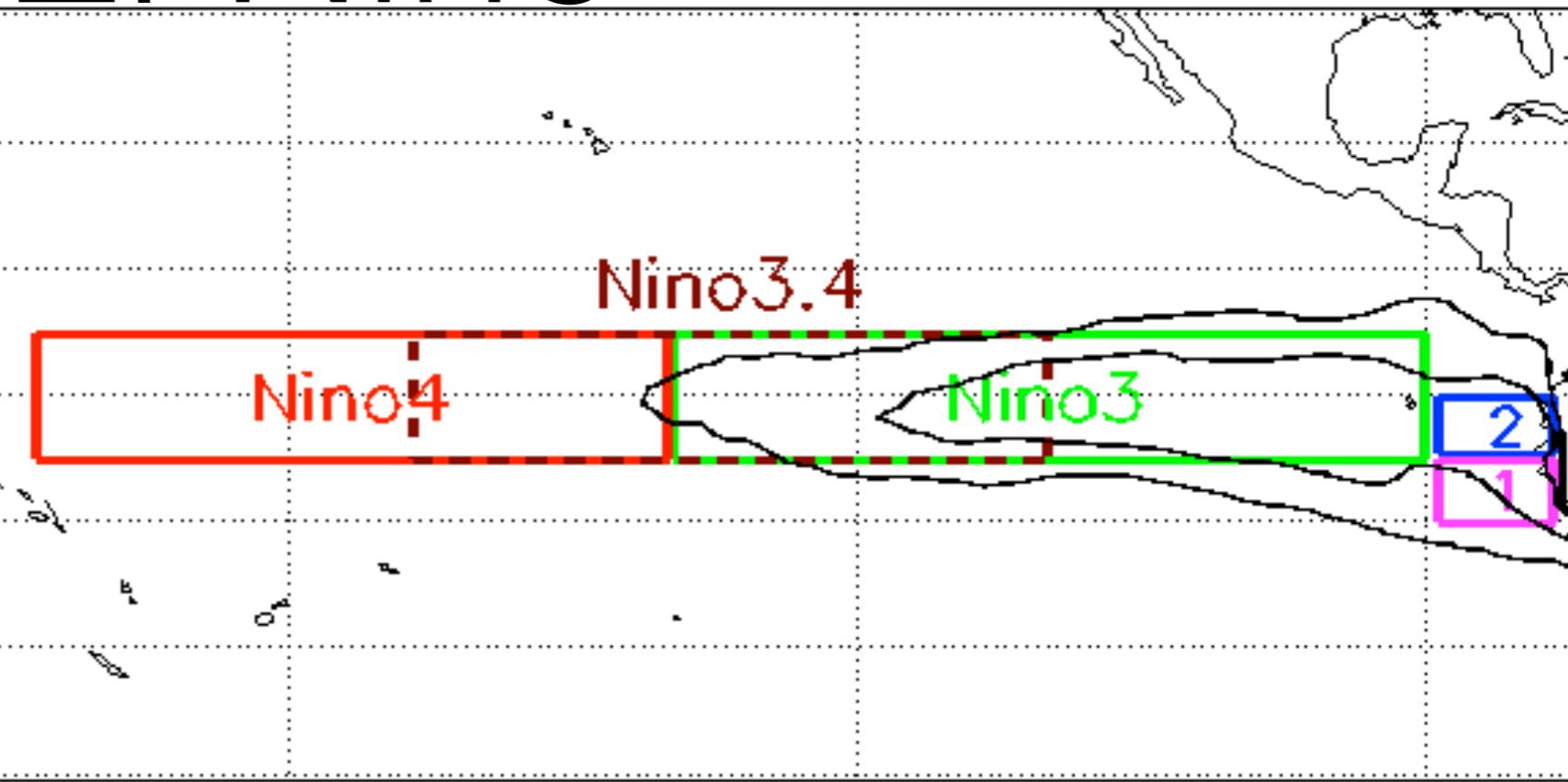


1 DEC 97

JPL



# El Nino





cow milk agriculture farm cattle livestock dairy  
beef hayfield field grass mammal pasture calf  
farmland rural animal pastoral bull grassland



cow beef agriculture cattle milk pasture mammal  
livestock farmland grass farm hayfield rural herd  
dairy pastoral grassland field calf bull



cow mammal pasture grass animal no person nature  
agriculture livestock hayfield cattle farm rural field  
milk grassland beef pastoral countryside



cow milk agriculture farm cattle livestock dairy  
beef hayfield field grass mammal pasture calf  
farmland rural animal pastoral bull grassland



cow beef agriculture cattle milk pasture mammal  
livestock farmland grass farm hayfield rural herd  
dairy pastoral grassland field calf bull



cow mammal pasture grass animal no person nature  
agriculture livestock hayfield cattle farm rural field  
milk grassland beef pastoral countryside



beach

sand

travel

no person

water

sea

seashore

depositphotos

depositphotos

depositphotos

depositphotos

depositphotos

depositphotos

depositphotos

depositphotos

depositphotos



no person

water

mammal

cattle

outdoors

cow

landscape

travel

sky

livestock



water

no person

beach

seashore

sea

sand

mammal

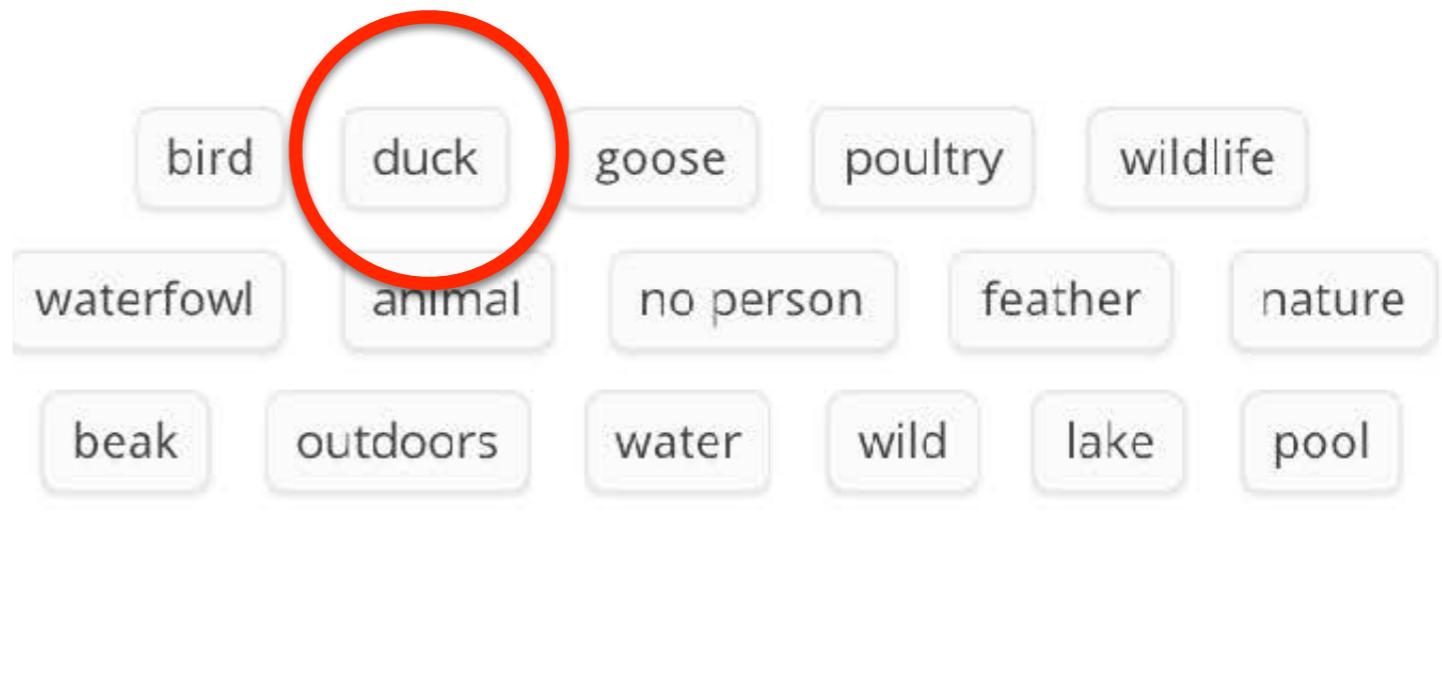
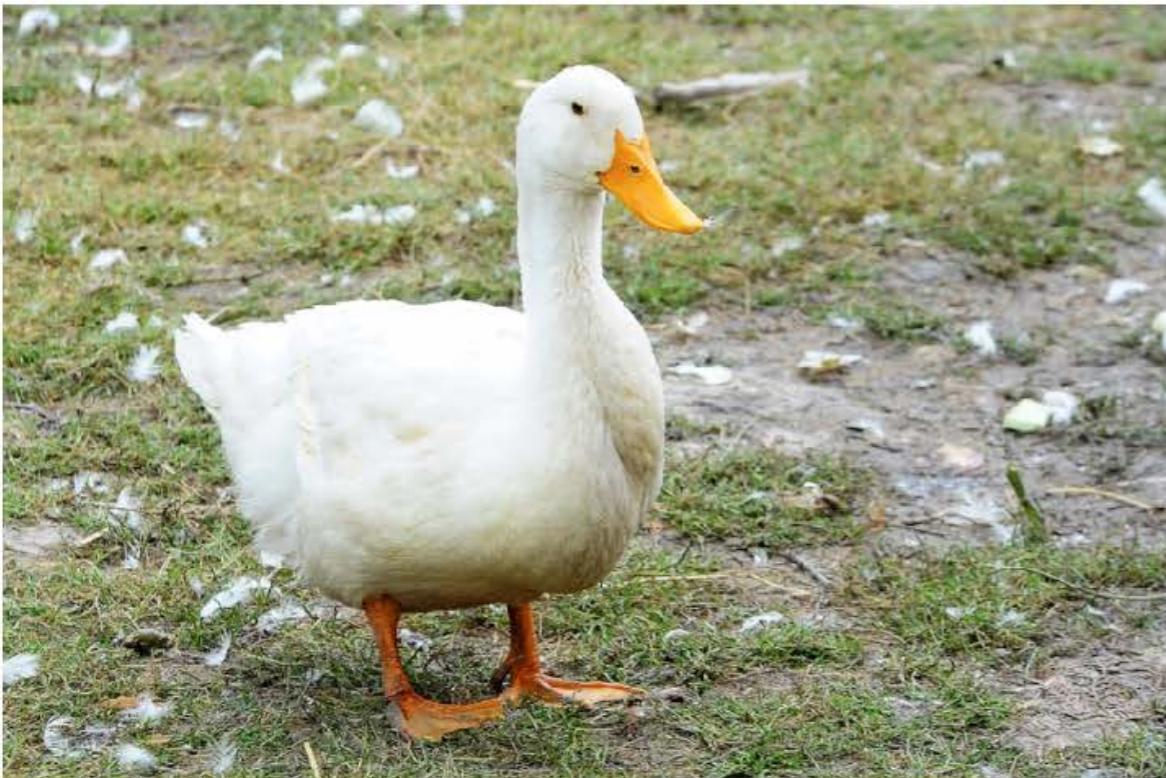
outdoors

travel

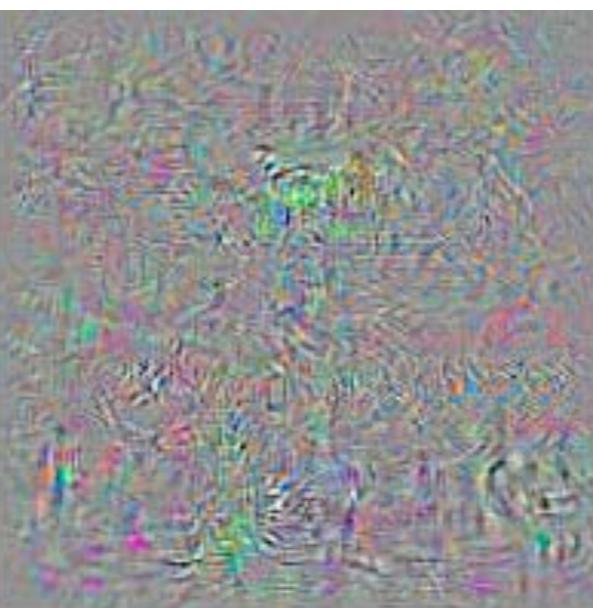
ocean

surf

sky

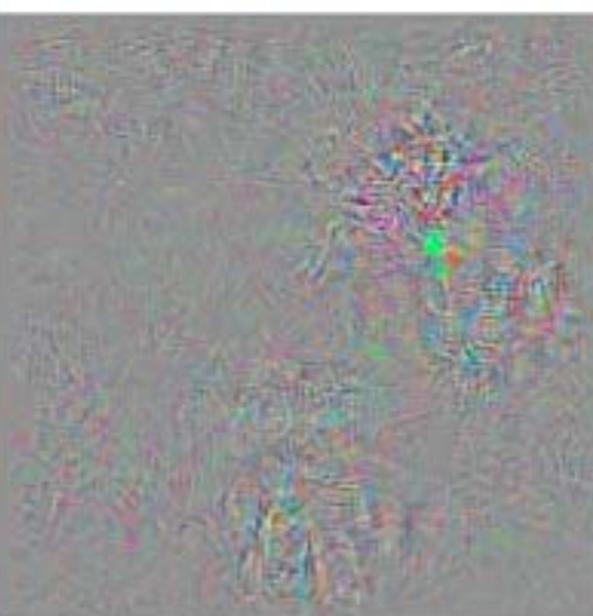
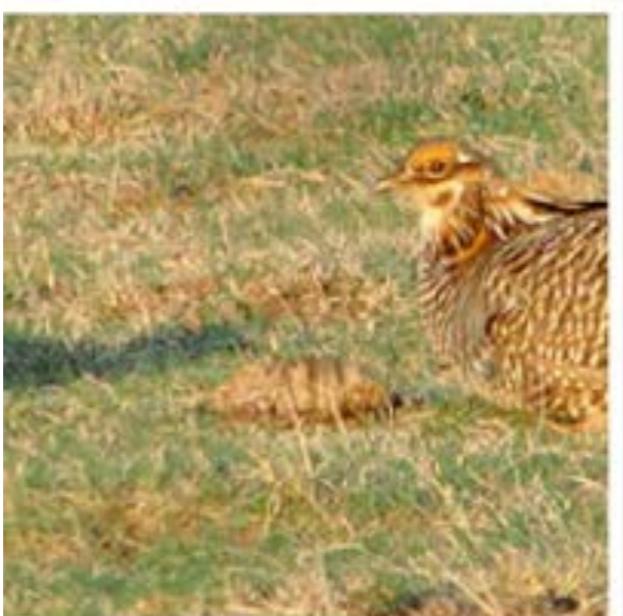


bird      no person      one



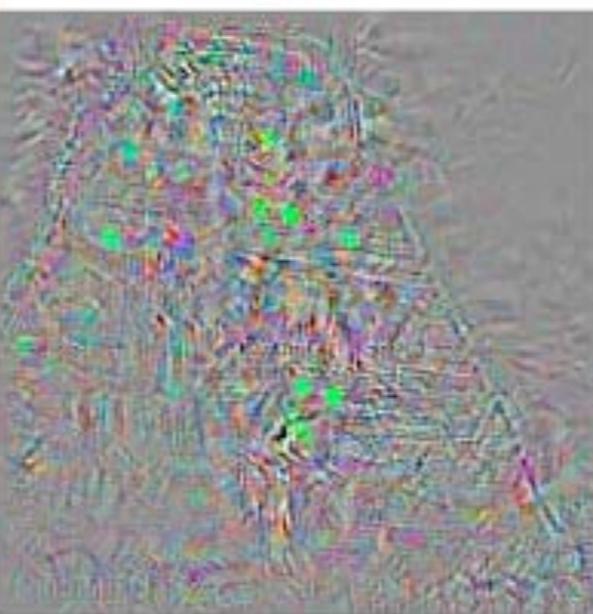
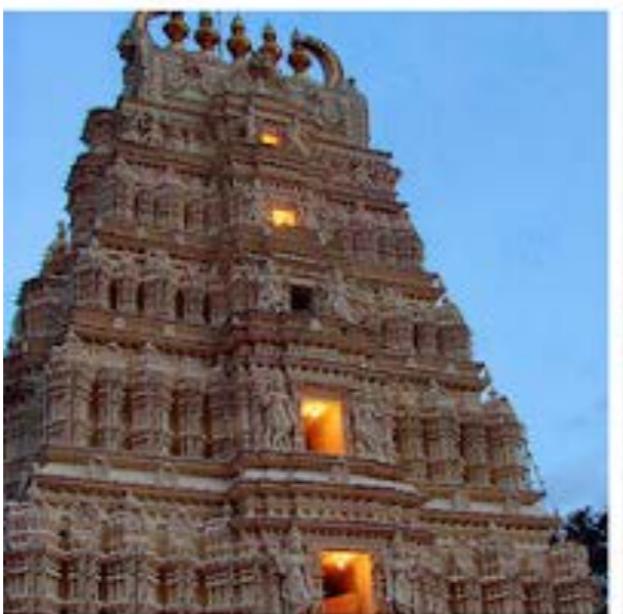
Schoolbus

Ostrich



Grouse

Ostrich



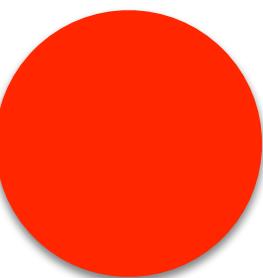
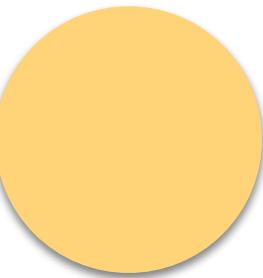
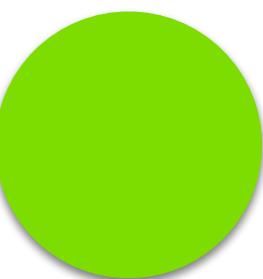
Pyramid

Ostrich

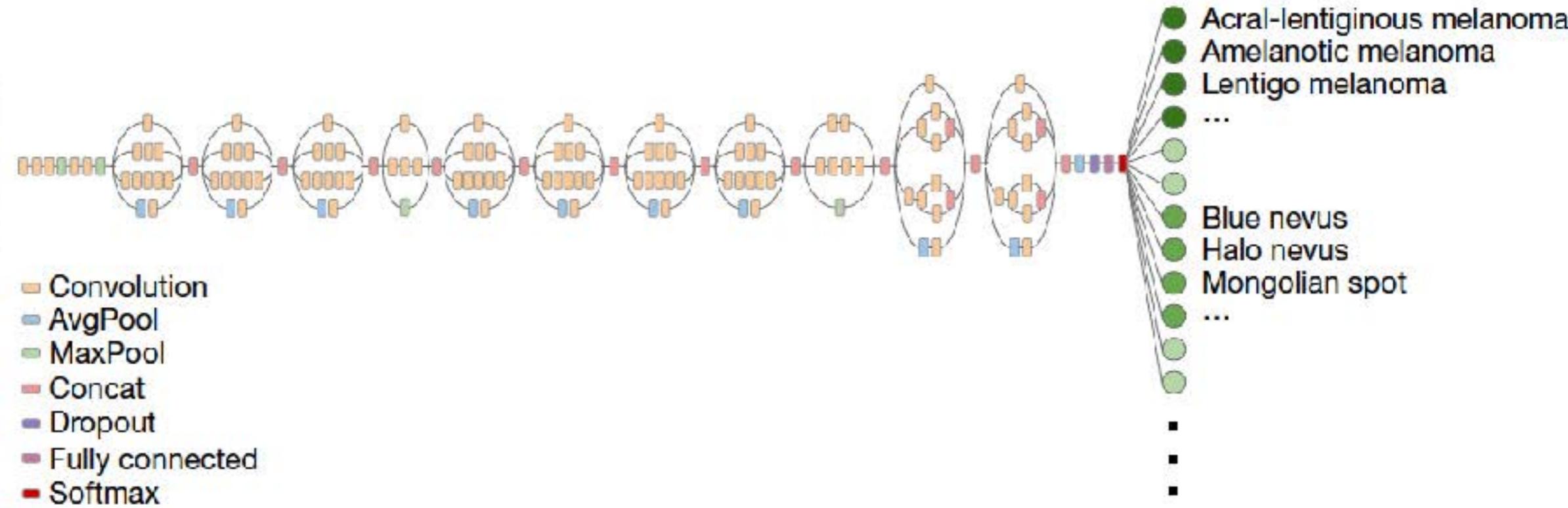
[Szegedy et al. 2014]

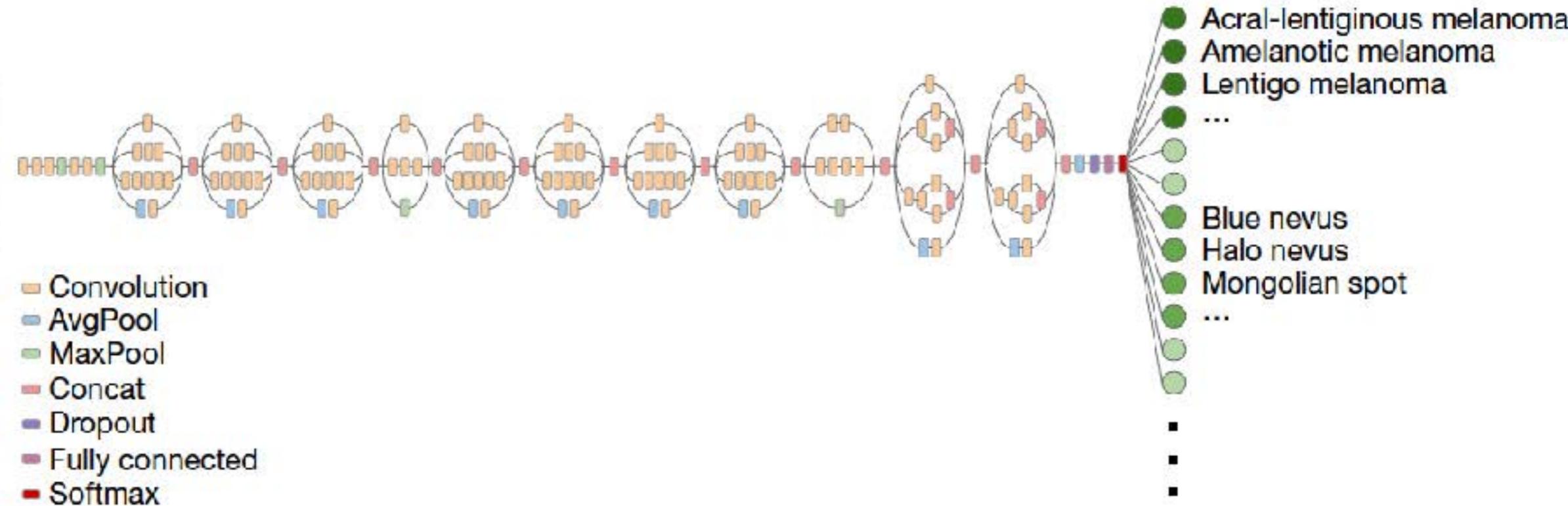
# Levels of understanding

- Memorization / recall
- Generalization / prediction
- Mechanisms / intervention



theory-driven design





## Questions:

Basic modules

Architecture

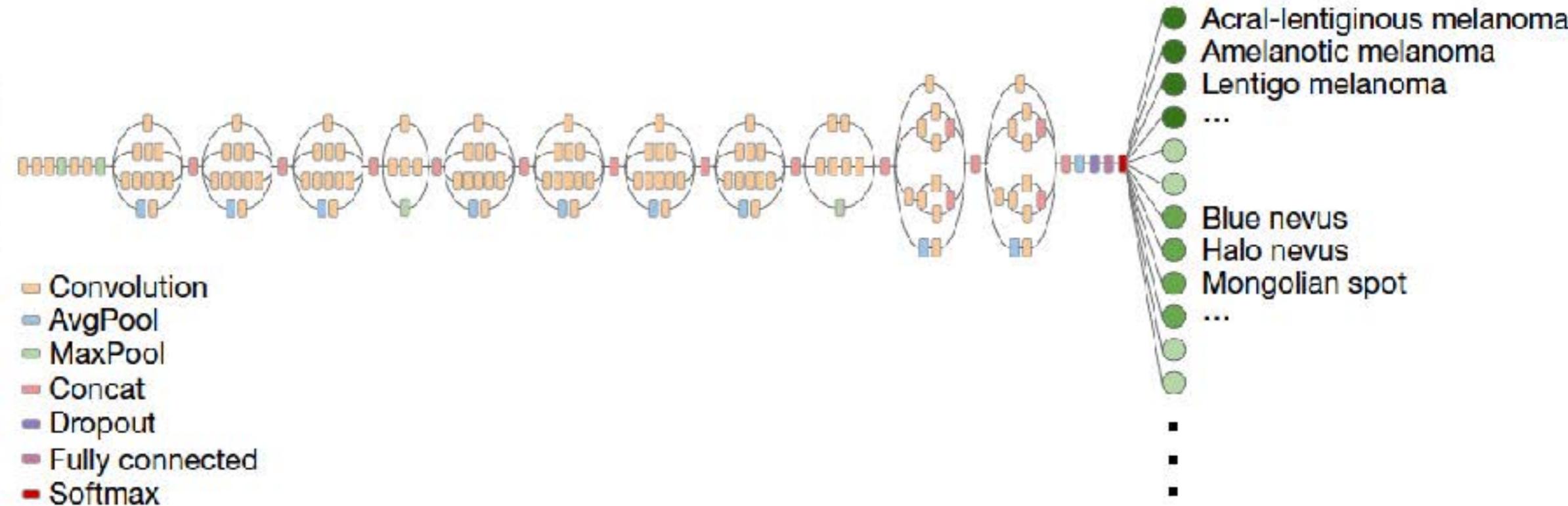
Optimization

$N \rightarrow 0$

Structure of data

Performance bounds

....



## Questions:

Basic modules

Architecture

Optimization

$N \rightarrow 0$

Structure of data

Performance bounds

Better  
understanding  
needed

# Conclusions

- Computer vision
- Learning-based approach
- Deep networks
- Practical results
- Open problems: N->0, causality, design
- Need theory